



INTERNATIONAL APPLICATION PUBLISHED UNDER THE PATENT COOPERATION TREATY (PCT)

<p><b>(51) International Patent Classification 5 :</b>  <b>G06F 15/42</b></p>	<p><b>A1</b></p>	<p><b>(11) International Publication Number:</b> <b>WO 94/11837</b></p> <p><b>(43) International Publication Date:</b> <b>26 May 1994 (26.05.94)</b></p>
<p><b>(21) International Application Number:</b> <b>PCT/US93/10507</b></p> <p><b>(22) International Filing Date:</b> <b>8 November 1993 (08.11.93)</b></p> <p><b>(30) Priority data:</b>  <b>07/975,526</b>      <b>12 November 1992 (12.11.92) US</b></p> <p><b>(71) Applicant (for all designated States except US):</b> <b>HITACHI CHEMICAL COMPANY, LTD. [JP/JP]; 1-1, 2-chome, Nishinshinjuku, Shinjuku-ku, Tokyo (JP).</b></p> <p><b>(72) Inventors; and</b>  <b>(75) Inventors/Applicants (for US only):</b> <b>MITSUHASHI, Masato [JP/US]; 8 Brookmont, Irvine, CA 92714 (US). COOPER, Allan [US/US]; 3607 North West 172nd Avenue, Bellevue, WA 98008 (US). WATERMAN, Michael [US/US]; 4336 Mentone Avenue, Culver City, CA 90232 (US). PEVZNER, Pavel [US/US]; 665 Cricklewood Drive, State College, PA 16803 (US).</b></p>	<p><b>(74) Agent:</b> <b>WAGNER, John, E.; Wagner &amp; Middlebrook, 3541 Ocean View Boulevard, Glendale, CA 91208 (US).</b></p> <p><b>(81) Designated States:</b> <b>CA, DE, GB, JP, KR, US, European patent (AT, BE, CH, DE, DK, ES, FR, GB, GR, IE, IT, LU, MC, NL, PT, SE).</b></p> <p><b>Published</b>  <i>With international search report.</i>  <i>With amended claims.</i></p>	
<p><b>(54) Title:</b> <b>OLIGOPROBE DESIGNSTATIONS: A COMPUTERIZED METHOD FOR DESIGNING OPTIMAL OLIGONUCLEOTIDE PROBES AND PRIMERS</b></p>		
<p><b>(57) Abstract</b></p> <p>There is disclosed herein an invention which relates to the fields of genetic engineering, microbiology, and computer science, that allows a user, whether a molecular biologist or a clinical diagnostician, to calculate and design extremely specific oligonucleotide sequences for DNA and mRNA hybridization procedures. The sequences designed with this invention may be used for medical diagnostic kits, DNA identification, and potentially continuous monitoring of metabolic processes in human beings. The key features design oligonucleotide sequences based on the GenBank database of DNA and mRNA sequences and examine candidate sequences for specificity or commonality with respect to a user-selected experimental preparation. Two models are available: a Mismatch Model, that employs hashing and continuous seed filtration, and an H-site Model, that analyzes candidate sequences for their binding specificity relative to some known set of mRNA or DNA sequences. The preferred embodiment of this computerized design tool is written in the Borland<sup>®</sup> C++ language and runs under Microsoft<sup>®</sup> Windows<sup>™</sup> on IBM<sup>®</sup> compatible personal computers.</p>		

**FOR THE PURPOSES OF INFORMATION ONLY**

Codes used to identify States party to the PCT on the front pages of pamphlets publishing international applications under the PCT.

AT	Austria	GB	United Kingdom	MR	Mauritania
AU	Australia	GE	Georgia	MW	Malawi
BB	Barbados	CN	Guinea	NE	Niger
BE	Belgium	GR	Greece	NL	Netherlands
BF	Burkina Faso	HU	Hungary	NO	Norway
BG	Bulgaria	IE	Ireland	NZ	New Zealand
BJ	Benin	IT	Italy	PL	Poland
BR	Brazil	JP	Japan	PT	Portugal
BY	Belarus	KE	Kenya	RO	Romania
CA	Canada	KG	Kyrgyzstan	RU	Russian Federation
CF	Central African Republic	KP	Democratic People's Republic of Korea	SD	Sudan
CG	Congo	KR	Republic of Korea	SE	Sweden
CH	Switzerland	KZ	Kazakhstan	SI	Slovenia
CI	Côte d'Ivoire	LI	Licthenstein	SK	Slovakia
CM	Cameroon	LK	Sri Lanka	SN	Senegal
CN	China	LU	Luxembourg	TD	Chad
CS	Czechoslovakia	LV	Latvia	TC	Togo
CZ	Czech Republic	MC	Monaco	TJ	Tajikistan
DE	Germany	MD	Republic of Moldova	TT	Trinidad and Tobago
DK	Denmark	MG	Madagascar	UA	Ukraine
ES	Spain	ML	Mali	US	United States of America
FI	Finland	MN	Mongolia	UZ	Uzbekistan
FR	France			VN	Viet Nam
GA	Gabon				

**OLIGOPROBE DESIGNSTATION: A COMPUTERIZED METHOD  
FOR DESIGNING OPTIMAL OLIGONUCLEOTIDE PROBES AND PRIMERS**

A portion of the disclosure of this patent document contains material which is subject to copyright protection. The copyright owner has no objection to the facsimile reproduction by anyone of the patent disclosure, as it appears in the Patent and Trademark Office patent files or records, but otherwise reserves all copyright rights whatsoever.

**BACKGROUND OF THE INVENTION**

This invention relates to the fields of genetic engineering, microbiology, and computer science, and more specifically to an invention that helps the user, whether they be a molecular biologist or a clinical diagnostician, to calculate and design extremely accurate oligonucleotide sequences for use as probes, for example for DNA and mRNA hybridization procedures, or as primers, for example for DNA amplification and extension using the polymerase chain reaction (PCR). In the following description, the design of probes has been discussed.

The oligonucleotide probes designed with this invention may be used to test for the presence of precursors of specific proteins in living tissues, or may be used for medical diagnostic kits, DNA identification, and potentially continuous monitoring of metabolic processes in human beings. The present implementation of this computerized design tool runs under Microsoft® Windows™ v. 3.1 (made by Microsoft Corporation of Redmond, Washington) on IBM® compatible personal computers (PC's).

Unless defined otherwise, all technical and scientific terms used herein have the same meaning as commonly understood by one of ordinary skill in the art to which this invention belongs. Although any methods and materials similar or equivalent to those described herein can be used in the practice or testing of the present invention, the preferred methods and materials are now described. All publications mentioned hereunder are incorporated herein by reference.

To isolate a specific gene for any particular purpose, a researcher first has to have some idea of what he or she is looking for. To do this, the researcher needs to have a probe, which acts like a molecular hook that can identify and latch onto (i.e., bind to or hybridize with) the desired gene in a crowd of many other genes. A researcher who can obtain an entire strand of mRNA can eventually find the gene from which it was copied, using complementary DNA (cDNA, which is a cloned equivalent

2

to RNA and somewhat equivalent to mRNA) as a probe to search through the great mass of genetic material and locate the desired original gene. cDNA essentially is manufactured or non-naturally occurring DNA from which all of the nonessential DNA has been removed. cDNA allows the researcher to concentrate entirely on the important portions of the gene being examined. The nonessential DNA regions are easy to recognize because when the gene is translated into protein, these regions do not wind up reflected in the protein sequence. These regions are called introns, or intervening regions. mRNA has no introns because they have been "spliced" out of the mRNA before translation. Thus, mRNA and cDNA contain only the essential information from a gene (called the exons). cDNA is the equivalent of mRNA with a complementary sequence, only the exons are present. cDNA may be produced by reverse transcription of mRNA.

The procedure of using cDNA from known mRNA as a probe to search through genetic material and locate the original gene is called molecular hybridization, and is currently one method of identifying specific genes. However, this method is less than perfect, can be extremely time consuming, and often is not even feasible because the researcher actually has to have an entire strand of cDNA from the desired gene before he or she can attempt to use this cDNA to locate and identify the particular gene. Thus, it is something of a circular problem. If the researcher cannot obtain an entire strand of mRNA or cDNA from the desired gene, then he or she must somehow design a probe from scratch to be used to identify that gene.

Oligonucleotide probes (that is, probes made up of a small number of nucleotides, such as 17 to 100), are increasingly being used to identify specific genes from genomic or cDNA libraries when the partial amino acid sequences is known. (von Heijne 1987, Ref. 15). This is a second method of determining a proper probe. Although the present implementation of this invention does not deal with cases in which the proteins have been sequenced, but rather only the DNA or mRNA, it is possible that this invention or a future implementation of it might be used with protein sequences. Such probes can also be used as primers which, when annealed to mRNAs, can be selectively extended into cDNAs. (von Heijne 1987, Ref. 15).

Because of these situations, the problem that the researcher faces is to discover or design a probe or mixture of probes that maximizes the researchers chances of successful hybridization while at the same time minimizing the amount of time and money that has to be spent on discovering or designing the probes. (von Heijne 1987,

3

Ref. 15). Researchers in the field have determined that computer analysis can greatly expedite and simplify the search for optimal probe sequences. (von Heijne 1987, Ref. 15). However, all of the search strategies known to the present inventors are time consuming (both CPU and user time) and may be somewhat inaccurate. As stated in von Heijne, "a true optimization of the probe in terms not only of degeneracy but in terms of length, codon usage, Guanine-Cytosine (GC) avoidance, and expected signal-to-noise ratio (hybridization to target over background) is a fairly complex problem, however, and does not seem to have been automated so far." (von Heijne 1987, Ref. 15). Various search strategies known and used in the field to identify and design probes are outlined in the following sources: Lewis (1986, Ref. 9), Raupach (1984, Ref. 11), Yang et al. (1984, Ref. 16), and Martin and Castro (1984, Ref. 10).

In the simplest version of a protein-related search strategy, the search procedure is limited to finding a set of probes of given lengths with the least possible degeneracy simply by scanning the amino acid sequence and noting the number of alternative codons in the corresponding oligonucleotide as the scan moves along the chain of nucleotides. (Lewis 1986). The researcher can also include codon usage statistics (because more than one codon can translate to the same amino acid), which would attach a probability-of-occurrence value to each probe. (Raupach 1984, Ref. 11).

A more advanced algorithm would allow the researcher to specify the way in which he or she plans to synthesize the probes (for example, by adding monomers or mixtures of monomers). It would also be easy for a researcher to add a rough estimate of the disassociation (or melting) temperatures of each probe to a program such as this.

One way to solve the problem of finding local similarities between two proteins being compared that has been discussed in the relevant literature is to use list-sorting or hashing routines. (von Heijne 1987, Ref. 15). These routines are based on the construction of a list or lookup table of k-letter words or k-tuples (i.e., all possible di- or trinucleotides), and the positions where they appear in the sequences being compared. This method is employed in some of the most extensively used "fast search" programs (see examples identified in von Heijne 1987, Ref. 15).

Two general methods of designing probes are common in the field, depending upon whether the researcher is trying to design a common probe or a specific probe. Common probes attempt to find common or consensus sequences among various species and among family genes. The first step in designing such a probe is to find the genes of interest. This may be done by performing a keyword or homology search against the

4

GenBank (a genome database available from IntelliGenics of Mountain View, CA) or a keyword search against MEDLINE (the database currently available from the U.S. National Library of Medicine under the data access system known as Dialog of Dialog Information Service, Inc., Palo Alto, CA) or by performing a homology analysis between one of the genes of interest and whole GenBank sequences. The next step is to retrieve all of the relevant genes of interest. In the third step, multiple alignment analysis can be done using a commercially available software package such as DNASIS (from Hitachi Software of Brisbane, California), which is an autoconnect program. In this step, the computer identifies which nucleotides are common among the requested sequences:

```

A1  A G G C C T C G G T T A G T T G G C C G T T G C C G A A A AA
    : : : : : : : : : : : : : : : : : : : : : :
A2  A G G C G T C G G T T A T T T G G G C C T T C C C A A T G TG
    : : : : : : : : : : : : : : : : : : : : : :
A3  A G G C G T C G G T T C T G T G G A A C T T C C C G A G G AA
    * * * * * * * * * * * * * * * * * * * *

```

\* = common among A1, A2, and A3

Alternatively, after homology analyses between two sequences are carried out, data from the multiple homology analyses can be combined. The researcher then manually has to find the common or consensus region:

```

A1  A G G C C T C G G T T A G T T G G C C G T T G C C G A A A AA
    : : : : : : : : : : : : : : : : : : : : : :
A2  A G G C G T C G G T T A T T T G G G C C T T C C C A A T G TG
    : : : : : : : : : : : : : : : : : : : : : :
A2  A G G C G T C G G T T A T T T G G G C C T T C C C A A T G TG
    : : : : : : : : : : : : : : : : : : : : : :
A3  A G G C G T C G G T T C T G T G G A A C T T C C C G A G G AA
    * * * * * * * * * * * * * * * * * * * *

```

\* = common among A1, A2, and A3

Next, the researcher would input the sequence of the common region into the program and then analyze the secondary structure (i.e., the stacking site and the hairpin structure). After this, the researcher manually would select several candidate probes (from five to ten) which contain the minimal hairpin structure and specific length according to the user's interest. A hairpin is an area in which a probe has "folded back" and one portion of the probe has hybridized with another portion of the same probe. The researcher would then perform a homology analysis between each candidate probe and all sequences in the GenBank to find all possible cross-hybridizable genes. Lastly,

5

the researcher manually would decide which is the best candidate probe by determining which probe is highly homologous among the group of interest, but quite different from other unrelated sequences in the GenBank.

The conventional methods for designing common oligonucleotide probes using currently available computer software have at least five problems: (1) they involve time consuming multiple processes; (2) it is difficult to control a significant variable, the melting temperature  $T_m$  of the oligonucleotide probes; (3) the methods do not recognize exons and introns and differentiate (thereby making it possible to have a designed probe that is identical to unrelated mRNA sequences); (4) the methods may miss short pieces of identical sequences; and (5) it is difficult to recognize multiple pieces of identical sequences in the gene.

The second method of designing probes that is common in the field involves designing specific probes. Specific probes attempt to find unique sequences among various species and among family genes and among published sequences in the GenBank. A specific probe is a probe that hybridizes with only one particular gene, thereby identifying the presence of that gene for the researcher. The procedure involves first finding the genes of interest (by performing a keyword search against the GenBank or against MEDLINE) and then retrieving all of the relevant genes of interest. A manual homology analysis between the gene of interest and whole sequences in the GenBank can be performed to find common and unique regions.

A1	A	G	G	C	C	T	C	G	G	T	T	A	G	T	T	G	G	C	C	G	T	T	G	C	C	G	A	A	A	A					
	:	:	:	:	:	:	:	:	:	:	:	:	:	:	:	:	:	:	:	:	:	:	:	:	:	:	:	:	:	:					
B1	A	G	G	C	G	T	C	G	G	T	T	A	T	T	G	T	G	G	T	C	T	C	C	C	C	A	A	T	G	TG					
	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-						
	common													* * * * *														unique							

Next, the researcher would input the sequence of the unique region into the program and then analyze the secondary structure. After this, the researcher would manually select several candidate probes which contain the minimal hairpin structure and specific length according to the user's interest. The researcher would then perform a homology analysis between each candidate probe and all sequences in the GenBank to find all possible cross-hybridizable genes. Lastly, the researcher manually would decide which is the best candidate probe by determining which probe does not have identical sequences in unrelated sequences in the GenBank.

All of the conventional methods for designing specific oligonucleotide probes known to the inventors using currently available computer software have at least four problems: (1) they involve time consuming multiple processes; (2) it is difficult to control the melting temperature  $T_m$  of the oligonucleotide probes; (3) the methods do not allow for quantification of uniqueness; and (4) there is no guarantee that the method will design the best possible probe.

None of the methods discussed in the literature discloses a system that may be used to design both common probes and extremely specific probes, especially a method that minimizes user and CPU time and is exceptionally accurate.

Programs currently used for rapid database similarity searches use either hashing strategies or statistical strategies. The hashing strategy is now being used for the detection of relatively short regions of similarity, while the statistical strategy is now being used for the detection of weaker and longer similarity regions. The Mismatch Model of this invention can be used for very strong similarity searches with running times faster than current hashing strategies.

The basic technologies behind the Mismatch Model used in this invention are hashing and continuous seed filtration, each general technology being known in the public domain and having been previously applied separately to non-genetic applications. To the best of the inventors' knowledge, these methods, used together, have never been suggested in other studies on optimal probe selection. The inventors' methods have a program performance of tens of seconds (CPU + I/O time) with a 1000 nucleotide query and all mammalian DNA on a SPARC station, and are even faster on the more common personal computer proposed herein.

The H-Site Model of this invention likewise is unique in that it offers a multitude of information on selected probes and original and distinctive means of visualizing, analyzing and selecting among candidate probes designed with the invention. Candidate probes are analyzed using the H-Site Model for their binding specificity relative to some known set of mRNA or DNA sequences, collected in a database such as the GenBank database. The first step involves selection of candidate probes at some or all the positions along a given target. Next, a melting temperature model is selected, and an accounting is made of how many false hybridizations each candidate probe will produce and what the melting temperature of each will be. Lastly, the results are presented to the researcher along with a unique set of tools for visualizing, analyzing and selecting among the candidate probes.



This invention is both much faster and much more accurate than the methods that are currently in use. It is unique because it is the only method that can find not only the most specific and unique sequence, but also the common sequences. Further, it allows the user to perform many types of analysis on the candidate probes, in addition to comparing those probes in various ways to the target sequences and to each other.

Therefore, it is the object of this invention to provide a practical and user-friendly system that will allow a researcher to design both specific and common oligonucleotide probes, and to do this in less time and with much more accuracy than currently done. For example, the current version of the GenBank contains over ninety (90) million nucleotides. It is thought that the human genome alone consists of three billion base pairs, and scientists have so far managed to decode the base sequence of only about 500 human genes, less than one percent of the total. Currently available searching strategies are limited in how many of the GenBank's sequences can be accessed and successfully searched, and how convenient and feasible such a search would be (in terms of both computer processor and human user time). It is also an object of this invention to allow the user to be able to run the program on more readily available and far less expensive computer hardware (i.e., a PC rather than a mainframe). This invention will remove those limits and allow genetic research to take a giant leap forward.

These and other advantages and objects of this invention will become apparent from the following detailed descriptions, drawings, and appended claims.

8

### BRIEF DESCRIPTION OF THE INVENTION

There is disclosed herein a system which allows the user to calculate and design extremely accurate oligonucleotide probes for DNA and mRNA hybridization procedures. The invention runs under Microsoft® Windows on IBM® compatible personal computers (PC's). Its key features design oligonucleotide probes based on the GenBank database of DNA and mRNA sequences and examine probes for specificity or commonality with respect to a user-selected experimental preparation of gene sequences. Hybridization strength between a probe and a subsequence of DNA or mRNA can be estimated through a hybridization strength model. Quantitatively, hybridization strength is given as the melting temperature  $T_m$ . Currently, two hybridization strength models are supported by this invention: 1) the Mismatch Model and 2) the H-Site Model. The user is allowed to select from the following calculations for each probe, results of which are available for display and analysis: 1) Sequence, Melting Temperature ( $T_m$ ) and Hairpin characteristics; 2) Hybridization with other species within the preparation mixture; and (3) Location and  $T_m$  for the strongest hybridizations. The results of the invention's calculations are then displayed on the Mitsuhashi Probe Selection Diagram (MPSD), which is a graphic display of all of the hybridizations of probes for the target mRNA with all sequences in the preparation.

The Main Dialog Window of the present implementation of this invention controls all user-definable settings. The user is offered a number of options at this window. The File option allows the user to print, print in color, save selected probes, and exit the program. The Preparation option allows the user to open and create preparation (PRP) files. The Models option allows the user to choose between the two hybridization models currently supported by the invention: 1) the H-Site Model and 2) the Mismatch Model. If the user selects the H-Site Model option, the user normally sets the following model parameters: 1) the melting temperature  $T_m$  for which probes are being designed (i.e., the melting temperature that corresponds to a particular experiment or condition the user desires to simulate); and 2) the nucleation threshold, which is the number of base pairs constituting a nucleation site. If the user selects the Mismatch Model option, the user normally sets the following model parameters: 1) probe length, which is the number of bases in probes to be considered; and 2) mismatch N, which is the maximum number of mismatches constituting a hybridization.

The Mismatch Model program is used to design DNA and mRNA probes, utilizing sequence database information from sources such as GenBank and other

databases with similar file formats. In the Mismatch Model, hybridization strength is related only to the number of base pair mismatches between a probe and its binding site. Generally, the more mismatches a user allows, the more probes will be found. The Mismatch Model does not take into account the Guanine-Cytosine (GC) content of candidate probes, as does the H-Site Model, discussed below, so there is no reflection or indication of the probe's binding strength. The basic technologies employed by this model are hashing and continuous seed filtration. Hashing involves the application of an algorithm or process to the records in a set of data to obtain a symmetric grouping of the records. When using an indexed set of data, hashing is the process of transforming a record key to an index value for storing and retrieving a record. Rosenberg (1984, Ref. 12)). The concept of continuous seed filtration is discussed in detail below.

The essence of the Mismatch Model is a fast process for doing exact and inexact matching between DNA and mRNA sequences to support the Mitsushashi Probe Selection Diagram (MPSD) and other types of analysis discussed above. The process used by the Mismatch Model is the Waterman-Pevzner Algorithm (the WPALG, which is named for two of the inventors), which is a computer-based probe selection process. Essentially, this is a combination of new and improved pattern matching processes. See Hume and Sunday (1991, Ref. 4), Landau *et al* (1986-1990, Refs. 6, 7, 8), Grossi and Luccio (1989, Ref. 3), and Ukkonen (1982, Ref. 14).

There are three principal programs that make up the Mismatch Model in this implementation of the invention. The first is designated by the inventors as "k\_diff." WPALG is used in k\_diff to find all locations of matches of length greater than or equal to one (1) (length is user-specified) with less than or equal to k number of mismatches (k is also user-specified) between the two sequences. If a candidate oligonucleotide probe fails to match that well, it is considered unique. k\_diff uses hashing and continuous seed filtration, and looks for homologs in GenBank and other databases with similar file formats. The technique of continuous seed filtration allows for much more efficient searching than previously implemented techniques. A seed is defined in this invention to be a subsequence of length equal to the longest exact match in the worst case scenario. For example, suppose the user selects a probe length (l) of 18, with 2 or fewer mismatches (k). If a match exists with 2 mismatches, then there must be a perfectly matching subsequence of length equal to 6. Once the seed length has been determined, the Mismatch Model looks at all substrings of that seed length (in this

example, that seed length would be 6), finds the perfectly matched base pair subsequence of length equals 6, and then looks to see if this subsequence extends to a sequence of length equal to the user selected probe length (i.e., 20 in this example). If so, a candidate probe has been found that meets the user's criteria.

Where the seed size is large, the program allocates a relatively large amount of memory for the hash table. This invention has an option that allows memory allocation for GenBank entries just once at the beginning of the program, instead of reallocating memory for each GenBank entry. This reduces input time for GenBank entries by as much as a factor of two (2), but the user needs to know the maximum GenBank entry size in advance to do this.

A probe is defined to hybridize if it has k or fewer mismatches in comparison with a target sequence from the database or file searched. Otherwise, it is non-hybridizing. The hit extension time for all appropriate parameters of the Mismatch Model has been found by experimentation to be less than thirty-five (35) seconds, except in one case where the minimum probe length (l) was set to 24 and the maximum number of mismatches (k) was set to four (4), which is a situation that is never used in real gene localization experiments because the hybridization conditions are too weak.

In this invention, the second hybridization strength model is termed the H-Site Model. One aspect of the H-Site Model uses a generalization of an experimental formula in general usage. The basic formula on which this aspect of the model is built is as follows:

$$T_m = \frac{81.5 - 16.6(\log[Na]) - .63 \%(\text{formamide}) + .41 (\%(\text{G} + \text{C})) - 600}{N}$$

In this formula, log[Na] is the sodium concentration, %(G + C) is the fraction of matched base pairs which are G-C complementary, and N is the probe length. In other words, this formula is an expression of the fact that melting temperature  $T_m$  is a function of both probe length and percent of Guanine-Cytosine (GC) content. This basic formula has been modified in this invention to account for the presence of mismatches. Each percent of mismatch reduces the melting temperature  $T_m$  by an average of 1.25 degrees (2 degrees C for an Adenine-Thymine mismatch, and 4 degrees C for a Guanine-Cytosine mismatch). This formula is, however, an approximation. The actual melting temperature might differ significantly from this approximation, especially for short probes or for probes with a relatively large number of mismatches.

//

Hybridization strength in the H-Site Model is related to each of the following factors: 1) "binding region"; 2) type of mismatch (GC or AT substitution); 3) length of the probe; 4) GC content of the binding region (since GC pairs have a stronger bond than AT pairs, thus requiring a higher melting temperature); and 5) existence of a "nucleation site" (an exactly matching subsequence). The type of mismatch and the GC content of the binding region each contribute to a candidate probe's binding strength, which can be compared to other candidate probes' binding strengths to enable the user to select the optimal probe.

The fundamental assumption of the H-Site Model is that binding strength is determined by a paired subsequence of the probe-species combination, called the binding region. If the binding region contains more GC pairs than AT pairs, the binding strength will be higher since the G and C bases (connected with three bonds) form a tighter bond than the A and T bases (connected with two bonds). Thus, G and C bases, and probes that are GC rich, require a higher melting temperature  $T_m$  and subsequently form a stronger bond. In the H-Site Model, and one of its unique features, the program designs optimal probes, ideally ones that do not have any mismatches, but if there are mismatches the H-Site Model takes these into account. With this model, a candidate probe can afford to have more mismatches involving the AT bases if there are more GC bases than AT bases in the probe. This is because this model looks primarily at regions of the candidate probe and target sequence that match and does not "penalize" the probe for areas that do not match. If the mismatches are located at either or both of the ends of the binding region, this has little effect. It is much more deleterious to have mismatches in the middle of the binding region, as this will significantly lower the binding strength of the probe.

The formula cited above for  $T_m$  applies within the binding region. The length of the probe is used to calculate percentages, but all other parameters of the formula are applied to the binding region only. The H-Site Model further assumes the existence of a nucleation site, which is a region of exact match. The length of this nucleation site may be set by the user. Typically, a value of 8 to 10 base pairs is used. To complete the H-Site Model, the binding region is chosen so as to maximize the melting temperature  $T_m$  among all regions containing a nucleation site, assuming one exists (otherwise,  $T_m=0$ ).

The H-Site Model is more complex than the Mismatch Model discussed above in that hybridization strength is modeled as a sum of signed contributions, with matches

12

generally providing positive binding energy and mismatches generally providing negative binding energy. The exact coefficients to be used depend only on the matched or mismatched pair. These coefficients may be specified by the user, although in the current version of this invention these coefficients are not explicitly user-selectable, but rather are selected to best fit the hybridization strength formulas developed by Itakura *et al* (1984, Ref. 5), Bolton and McCarthy (1962, Ref. 2), Benner *et al* (1973, Ref. 1), and Southern (1975, Ref. 13).

A unique aspect of the H-Site Model is that hybridization strength is defined to be determined by whatever the optimal binding region between the candidate probe and binding locus. This binding region is called the hybridization site, or h-site, and is selected so as to maximize overall hybridization strength, so that mismatches outside the binding region do not detract from the estimated hybridization strength. Several other unique features of the H-Site Model include the fact that it is more oriented toward RNA and especially cDNA sequences than DNA sequences, and the fact that the user has control over preparation and environmental variables. The first feature allows the user to concentrate on "meaningful" sequences, rather than having to sort through all of a DNA sequence (including the introns). The second feature allows the user to more accurately simulate laboratory conditions and more closely correspond with any experiments he or she is conducting. Further, this implementation of the invention does some preliminary preprocessing of the GenBank database to sort out and select the cDNA sequences. This is done by locating a keyword (in this case CDS) in each GenBank record, thereby eliminating any sequences containing introns.

The Mitsuhashi Probe Selection Diagram (MPSD), FIG. 4, is the third key feature of this invention, as it is a unique way of visualizing the results of the probe designing performed by the Mismatch and H-Site Models. It is a graphic display of all of the hybridizations of candidate oligonucleotide probes for the target mRNA with all sequences in the preparation. Given a gene sequence database and a target mRNA sequence, the MPSD graphically displays all of the candidate probes and their hybridization strengths with all sequences from the database. In the present implementation, each melting temperature  $T_m$  is displayed as a different color, from red (highest  $T_m$ ) to blue (lowest  $T_m$ ). The MPSD allows the user to see visually the number of false hybridizations at various temperatures for all candidate probes, and the sources of these false hybridizations (with a loci and sequence comparison). A locus

13

may be a specific site or place, or, in the genetic sense, a locus is any of the homologous parts of a pair of chromosomes that may be occupied by allelic genes.

**BRIEF DESCRIPTION OF THE DRAWING**

This invention may be more clearly understood from the following detailed description and by reference to the drawing in which:

FIG. 1 is a simplified block diagram of a computer system illustrating the overall design of this invention;

FIG. 2 is a display screen representation of the main dialog window of this invention;

FIG. 3 is a flow chart of the overall invention illustrating the program, and the invention's sequence and structure;

FIG. 4 is a display screen representation of the Mitsuhashi probe selection diagram;

FIG. 5 is a display screen representation of the probeinfo and matchinfo window;

FIG. 6 is a display screen representation of the probesedit window;

FIG. 6a is a printout of the probesedit output file;

FIG. 7 is a flow chart of the overall k\_diff program of the Mismatch Model of this invention, including its sequence and structure;

FIG. 8 is a flow chart of the k\_diff module of this invention;

FIG. 9 is a flow chart of the hashing module of this invention;

FIG. 10 is a flow chart of the tran module of this invention;

FIG. 11 is a flow chart of the let\_dig module of this invention;

FIG. 12 is a flow chart of the update module of this invention;

FIG. 13 is a flow chart of the assembly module of this invention;

FIG. 14 is a flow chart of the seqload module of this invention;

FIG. 15 is a flow chart of the read1 module of this invention;

FIG. 16 is a flow chart of the dig\_let module of this invention;

FIG. 17 is a flow chart of the q\_colour module of this invention;

FIG. 18 is a flow chart of the hit\_ext module of this invention;

FIG. 19 is a flow chart of the colour module of this invention;

FIG. 20 is a printout of a sample file containing the output of the Mismatch Model program of this invention;

FIG. 21 is a flow chart of the H-Site Model, stage I, covering the creation of a preprocessed preparation file of this invention;

FIG. 22 is a flow chart of the H-Site Model, stage II, covering the preparation of the target sequence(s);



15

FIG. 23 is a flow chart of the H-Site Model, stage III, covering the calculation of MPSD data;

FIG. 24a is a printout of a sample file containing output of the Mismatch Model program;

FIG. 24b is a printout of a sample file containing output of the H-Site Model program;

FIG. 25 is a flow chart of the processing used to create the Mitsuhashi probe selection diagram (MPSD);

FIG. 26 is a flow chart of processing used to create the matchinfo window;

FIG. 27 is a printout of a sample target species file;

FIG. 28 is a printout of a sample preparation file.

16

### DETAILED DESCRIPTION OF THE INVENTION

This invention is employed in the form best seen in FIG. 1. There, the combination of this invention consists of an IBM® compatible personal computer (PC), running software specific to this invention, and having access to a distributed database with the file formats found in the GenBank database and other related databases.

The preferred computer hardware capable of operating this invention involves of a system with at least the following specifications (FIG. 1): 1) an IBM® compatible PC, generally designated 1A, 1B, and 1C, with an 80486 coprocessor, running at 33 Mhz or faster; 2) 8 or more MB of RAM, 1A; 3) a hard disk 1B with at least 200 MB of storage space, but preferably 1 GB; 4) a VGA color monitor 1C with graphics capabilities of a size sufficient to display the invention's output in readable format, preferably with a resolution of 1024 x 768; and 5) a 580 MB CD ROM drive 5 (1B of FIG. 1 generally refers to the internal storage systems included in this PC, clockwise from upper right, two floppy drives, and a hard disk). Because the software of this invention preferably has a Microsoft® Windows™ interface, the user will also need a mouse 2, or some other type of pointing device.

The preferred embodiment of this invention would also include a laser printer 3 and/or a color plotter 4. The invention may also require a modem (which can be internal or external) if the user does not have access to the CD ROM versions of the GenBank database 8 (containing a variable number of gene sequences 6). If a modem is used, information and instructions are transmitted via telephone lines to and from the GenBank database 8. If a CD ROM drive 5 is used, the GenBank database (or specific portions of it) is stored on a number of CDs.

The computer system should have at least the Microsoft® DOS v. 5.0 operating system running Microsoft® Windows™ v. 3.1. All of the programs in the preferred embodiment of the invention are written in the Borland® C++ (made by Borland International, Inc., of Scotts Valley, CA) computer language. It must be recognized that subsequently developed computers, storage systems, and languages may be adapted to utilize this invention and vice versa.

This invention is designed to enable the user to access DNA, mRNA and cDNA sequences stored either in the GenBank or in databases with similar file formats. GenBank is a distributed flat file database made up of records, each record containing a variable number of fields in ASCII file format. The stored database itself is distributed, and there is no one database management system (DBMS) common to even a majority of its users. One general format, called the line type format, is used both for

17

the distributed database and for all of GenBank's internal record keeping. All data and system files and indexes for GenBank are kept in text files in this line type format.

The primary GenBank database is currently distributed in a multitude of files or divisions, each of which represents the genome of a particular species (or at least as much of it as is currently known and sequenced and publicly available). The GenBank provides a collection of nucleotide sequences as well as relevant bibliographic and biological annotation. Release 72.0 (6/92) of the GenBank CD distribution contains over 71,000 loci with a total of over ninety-two (92) million nucleotides. GenBank is distributed by IntelliGenetics, of Mountain View, CA, in cooperation with the National Center for Biotechnology Information, National Library of Medicine, in Bethesda, MD.

1. Overall Description of the Invention

a. General Theory

The intent of this invention is to provide one or more fast processes for performing exact and inexact matching between DNA sequences to support the Mitsuhashi Probe Selection Diagram (MPSD), discussed below, and other analysis with interactive graphical analysis tools. Hybridization strength between a candidate oligonucleotide probe and a subsequence of DNA, mRNA or cDNA can be estimated through a hybridization strength model. Quantitatively, hybridization strength is given as the melting temperature  $T_m$ . Currently, two hybridization strength models are supported by the invention: 1) the Mismatch Model and 2) the H-Site Model.

b. Inputs

i. Main Dialog Window

The Main Dialog Window, FIG. 2, controls all user-definable settings. This window has a menu bar offering five options: 1) File 10; 2) Preparation 20; 3) Models 30; 4) Experiment 40; and 5) Help 50. The File 10 option allows the user to print, print in color, save selected probes, and exit the program. The Preparation 20 option allows the user to open and create preparation (PRP) files.

The Models 30 option allows the user to choose between the two hybridization models currently supported by the invention: 1) the H-Site Model 21 and 2) the Mismatch Model 25. If the user selects the H-Site Model 21 option, the left hand menu of FIG. 2C is displayed and the user sets the following model parameters: 1) the melting temperature  $T_m$  22 for which probes are being designed (i.e., the melting temperature that corresponds to a particular experiment or condition the user desires

18

to simulate); and 2) the nucleation threshold 23, which is the number of base pairs constituting a nucleation site. If the user selects the Mismatch Model 25 option, the right hand menu of FIG. 2C is displayed and the user sets the following model parameters: 1) probe length 26, which is the number of base pairs in probes to be considered; and 2) mismatch N 27, which is the maximum number of mismatches constituting a hybridization. Computation of the user's request will take longer with the H-Site Model if the threshold 23 setting is decreased and with the Mismatch Model if the number of mismatches K 27 is increased.

In addition, for both Model options the user chooses the target species 11 DNA or mRNA for which probes are being designed and the preparation 12, a file of all sequences with which hybridizations are to be calculated. A sample of a target species file is shown in FIG. 27 (humbjnx.cds), while a sample of a preparation file is shown in FIG. 28 (junmix.seq). Each of these inputs is represented by a file name and extension in general DOS format. In the target species and preparation fields, the file format follows the GenBank format, and each of the fields includes a default file extension. Pressing the "OK" button 41 of FIG. 2C will cause the processing to begin, and pressing the "Cancel" button 43 will cause it to stop.

The Experiment 40 option and the Help 50 option are expansion options not yet available in the current implementation of the invention.

#### c. Processing

FIG. 3 is a flow chart of the overall program, illustrating its sequence and structure. Generally, the main or "control" program of the invention basically performs overall maintenance and control functions. This program, as illustrated in FIG. 3, accomplishes the general housekeeping functions 51, such as defining global variables. The user-friendly interface 53, carries out the user-input procedures 55, the file 57 or database 59 access procedures, calling of the model program 62 or 63 selected by the user, and the user-selected report 65 or display 67, 69, 71 and 73 features. Each of these features is discussed in more detail in later sections, with the exception of the input procedures, which involves capturing the user's set-up and control inputs.

#### d. Outputs

##### i. The Mitsuhashi Probe Selection Diagram Window

The Mitsuhashi Probe Selection Diagram (MPSD), FIG. 4, is a key feature of the invention as it is a unique way of visualizing the results of the program's calculations. It is a graphic display of all of the hybridizations of probes for the target mRNA with

19

all sequences in the preparation. In other words, given a sequence database and a target mRNA, the MPSD graphically displays all of the candidate probes and their hybridization strengths with all sequences from the sequence database. The MPSD allows the user to see visually the number of false hybridizations at various temperatures for all candidate probes, and the sources of these false hybridizations (with a loci and sequence comparison).

For each melting temperature  $T_m$  of interest, a graphical representation of the number of hybridizations for each probe is displayed. In the preferred embodiment, this representation is color coded. In this implementation of the invention, the color red 123 identifies the highest melting temperature  $T_m$  and the color blue 124 identifies the lowest melting temperature  $T_m$ . Each mismatch results in a reduction in  $T_m$ .  $T_m$  is also a function of probe length and percent content of GC bases. Within the window, the cursor 125 shape is changed from a vertical line bisecting the screen to a small rectangle when the user selects a particular probe. The current probe is defined to be that probe under the cursor position (whether it be a line or a rectangle) in the MPSD window. More detailed information about the current probe is given in the ProbeInfo and MatchInfo windows, discussed below. Clicking the mouse 2 once at the cursor 125 selects the current probe. Clicking the mouse 2 a second time deselects the current probe. Moving the cursor across the screen causes the display to change to reflect the candidate probe under the current cursor position.

The x-axis 110 of the MPSD, FIG. 4, shows the candidate probes' starting positions along the given mRNA sequence. The user may "slide" the display to the left or right in order to display other probe starting positions. The y-axis 115 of the MPSD displays the probe specificity, which is calculated by the program.

The menu options 116, 117, 118, 119, and 120 available to the user while in the MPSD, FIG. 4, are displayed along a menu bar at the top of the screen. The user can click the mouse 2 on the preferred option to briefly display the option choices, or can click and hold the mouse button on the option to allow an option to be selected. The user may also type a combination of keystrokes in order to display an option in accordance with well-known computer desk top interface operations. This combination usually involves holding down the ALT key while pressing the key representing the first letter of the desired option (i.e., F, P, M, E or H).

The File option 116 allows the user to specify input files and databases. The Preparation option 117 allows the user to create a preparation file summarizing the

20

sequence database. The Models option 118 allows the user to specify the hybridization model (i.e., H-Site or Mismatch) and its parameters. The Experiment option 119 and the Help option 120 are not available in the current implementation of this invention. These options are part of the original Main Dialog Window, FIG. 2.

Areas on the graphical display of the MPSD, FIG. 4, where the hybridizations for the optimal probes are displayed are lowest and most similar, such as shown at 121, indicate that the particular sequence displayed is common to all sequences. Areas on the graphical display of the MPSD where the hybridizations for the optimal probes are displayed are highest and most dissimilar, such as shown at 122, indicate that the particular sequence displayed is extremely specific to that particular gene fragment. The high points on the MPSD show many loci in the database, to which the candidate probe will hybridize (i.e., many false hybridizations). The low points show few hybridizations, at least relative to the given database. In other words, the sequence shown at 121 would reflect a probe common to all of the gene fragments tested, such that this probe could be used to detect each of these genes. The sequence shown at 122 would reflect a probe specific to the particular gene fragment, such that this probe could be used to detect this particular gene and no others.

ii. The ProbeInfo and MatchInfo Window

The combined ProbeInfo and MatchInfo Window, FIG. 5, displays detailed information about the current candidate probe. The upper portion of the window is the ProbeInfo window, and the lower portion is the MatchInfo window. The ProbeInfo window portion displays the following types of information: the target locus (i.e., the mRNA, cDNA, or DNA from which the user is looking for probes) is displayed at 131, while the preparation used for hybridizations is displayed at 132. In the example shown in FIG. 5, the target locus 131 is the file named HUMBJUNX.CDS, which is shown as being located on drive F in the subdirectory MILAN. The preparation 132 is shown as being the file designated JUNMIX.PRP, which is also shown as being located on drive F in the subdirectory MILAN. The JUNMIX.PRP preparation in this example is a mixture of human and mouse jun loci.

The current and optimal probe's starting position is shown at 135. The current candidate oligonucleotide probe is defined at 136, and is listed at 137 as having a length of 21 bases. The melting temperature for the probe 136 as hybridized with the targets is shown in column 140. The melting temperature for the optimal probe is given as 61.7 degrees C at 138. The ProbeInfo Window FIG. 5 also displays hairpin characteristics

2/

of the probe at 139. In the example shown, the ProbeInfo Window shows that there are four (4) base pairs involved in the worst hairpin, and that the worst hairpin has a length of one (1) (see FIG. 5, at 139).

The MatchInfo Window portion displays a list of hybridizations between the current probe and species within the preparation file, including hybridization loci and hybridization temperatures. The hybridizations are listed in descending order by melting temperature. The display shows the locus with which the hybridization occurs, the position within the locus, and the hybridization sequence.

In the MatchInfo window portion, the candidate probe 136 is shown at 150 as hybridizing completely with a high binding strength. This is because the target DNA is itself represented in the database in this case, so the candidate probe is seen at 150 to hybridize with itself (a perfect hybridization). The locus of each hybridization from the preparation 132 are displayed in column 141, while the starting position of each hybridization is given in column 142. The calculated hybridizations are shown at 145.

### iii. The ProbesEdit Window

The ProbesEdit Window, FIG. 6, is a text editing window provided for convenient editing and annotation of the invention's text file output. It is also used to accumulate probes selected from the MPSD, FIG. 4, by mouse 2 clicks. Standard text editing capabilities are available within the ProbesEdit Window. The user may accumulate selected probes in this window (see 155 for an example) and then save them to a file (which will bear the name of the preparation sequence with the file extension of "prb" 156, or may be another file name selected by the user). A sample of this file is shown in FIG. 6A.

### iv. Miscellaneous Output

The present embodiment of this invention also creates two output files, currently named "test.out" and "test1.out", depending upon which model the user has selected. The first file, "test.out", is created with both the Mismatch Model and the H-Site Model. This file is a textual representation of the Mitsuhashi Probe Selection Diagram (MPSD). It breaks the probe sequence down by position, length, delta Tm, screensN, and the actual probe sequence (i.e., nucleotides). An example of this file created by the Mismatch Model is shown in FIG. 20, and example created by the H-Site Model is shown in FIG. 24A. The second file, "test1.out", is created only by the H-Site Model. This file is a textual representation of the ProbeInfo and MatchInfo window that captures all hybridizations, along with their locus, starting position, melting temperature,

22

and possible other hybridizations. A partial example of this file is shown in FIG. 24B (10 pages out of a total of 190 pages created by the H-Site Model).

## 2. Description of the Mismatch Model Program

### a. Overview

In this invention, one of the hybridization strength models is termed the Mismatch Model (see FIG. 2 for selection of this model). The basic operation of this model involves the techniques of hashing and continuous seed filtration, as defined earlier and described in more detail below. The essence of the Mismatch Model is a fast process for doing exact and inexact matching between DNA and mRNA sequences to support the Mitsuhashi Probe Selection Diagram (MPSD). There are a number of modules in the present implementation of the Mismatch Model contained in this invention, the most significant of which are shown in the flow chart in FIG. 7 and in more detail in FIGS. 8 through 18. The main k\_diff module shown in the flow chart in FIG. 8 is a structured program that provides overall control of the Mismatch Model, calling various submodules that perform different functions.

### b. Inputs

The user-selected input variables for this model are minimum probe length 26 (which is generally from 18 to 30) and maximum number of mismatches 27 (which generally is from 1 to 5). These inputs are entered by the user in the Main Dialog Window, FIG. 2C.

### c. Processing

#### i. k\_diff Program

Some terms of art need to be defined before the processing performed by this module can be explained. A hash table basically is an array or table of data. A linked list is a classical data structure which is a chain of linked entries and involves pointers to other entry structures. Entries in a linked list do not have to be stored sequentially in memory, as is the case with elements contained in an array. Usually there is a pointer to the list associated with the list, which is often initially set to point to the start of the list. A pointer to a list is useful for sequencing through the entries in the list. A null pointer (i.e., a pointer with a value of zero) is used to mark the end of the list.

As the flow charts in FIGS. 7 and 8 illustrate, the general process steps and implemented functions of this model can be outlined as follows:



## 23

Step 1: First, create a hash table and linked list from the query (FIG. 7, hashing module 222).

Step 2: Next, while there are still GenBank entries available for searching (FIG. 7, assembly module 230):

Step 2a: Read the current GenBank entry (record) sequence of user-specified length (FIG. 7, seqload module 232), or read the current sequence (record) from the file selected by the user (FIG. 7, read1 module 234).

Step 2b: For the current sequence for each position of the sequence from the first position (or nucleotide) to the last position (or nucleotide) (incrementing the position number once each iteration of the loop) (FIG. 7, q\_colour module 242),

Step 2c: set the variable dna\_hash equal to the hash of the current position of the current sequence (FIG. 7, q\_colour module 242).

Step 2d: While not at the end of the linked list for dna\_hash (FIG. 7, q\_colour module 242),

Step 2e: set the query\_pos equal to the current position of dna\_hash in the linked list (FIG. 7, q\_colour module 242) and

Step 2f: Extend the hit with the coordinates (query\_pos, dna\_pos) (FIG. 7, hit\_ext module 244),

Step 2g: If there exists a k\_mismatch in the current extended hit (FIG. 7, colour module 246), then

Step 2h: print the current hit (FIG. 7, q\_colour module 242), and repeat from Step 2.

As this illustrates, there are three (3) basic looping or iteration processes with functions being performed based on variables such as whether the GenBank section end has been reached (the first "WHILE" loop, Step 2), whether the end of the current DNA entry has been reached (the "FOR" loop, Step 2b), and whether the end of the dna\_hash linked list has been reached (the second "WHILE" loop, Step 2d). A "hit" will only be printed if there are k\_mismatches in the current extended hit.

FIGS. 8 through 18 illustrate the functions of each of the modules of the present embodiment of this invention, all of which were generalized and summarized in the description above. FIG. 8, which outlines the main "k\_diff" module, shows that this

24

module is primarily a program organization and direction module, in addition to performing routine "housekeeping" functions, such as defining the variables and hash tables 251, checking if the user-selected gene sequence file is open 252, extracting needed identification information from the GenBank 253, and ensuring valid user input 254. This module also performs a one-time allocation of memory for the gene sequences, and allocates memory for hit information, hashing, hybridization and frequency length profiles and output displays, 255 & 256. The "k\_diff" module also initializes or "zeros out" the hashing table, the linked hashing list and the various other variables 257 in preparation for the hashing function. In addition, this module forms the hash tables 258 and extracts a sequence and finds the sequence length 259.

One of the most important functions performed by the "k\_diff" module is to define the seed (or kernel or k\_tuple) size. This is done by setting the variable k\_tuple equal to  $(\text{min\_probe\_length} - \text{max\_mismatch\_}) / (\text{max\_mismatch} + 1)$  FIG. 8 at 265. Next, if the remainder of the aforementioned process is not equal to zero 266, then the value of the variable k\_tuple is incremented by one 267. The resulting value is the size of the seed. The module then reads the query 268 and copies the LOCUS name 269 for identification purposes (a definition of the term locus is given earlier in the specification).

The "k\_diff" module FIG. 8 also calls the "assembly"-module 260, writes the results to a file 261a, plots the results 261b (discussed below), calculates the hairpin characteristics 262 (i.e., the number of base pairs and the length of the worst hairpin) and the melting temperature (Tm) for each candidate probe 263, and saves the results to a file 264.

The screen graphs are plotted 261b by converting the result values to pixels, filing a pixel array and performing a binary search into the pixel array. Next, given the number of pixels per probe position and which function is of interest to the user (i.e., the three mismatch match numbers), the program interpolates the values at the value of (pixelsPerPositionN-1) and computes the array of pixel values for drawing the graph. These values are then plotted on the MPSD.

The "hashing" module, FIG. 9, performs hashing of the query. In other words, it creates the hash table and linked list of query positions with the same hash. The variable has\_table[i] equals the position of the first occurrence of hash i in the query. If i does not appear in the query, hash\_table[i] is set to zero.

25

The "tran" module, FIG. 10, is called by the "hashing" module 271, and performs the hashing of the sequence of  $k\_tuple$  (kernel or seed) size. If the  $k\_tuple$  exists (i.e., its length is greater than zero), the variable  $uns$  is set equal to  $uns*ALF+p$  291. The variable  $p$  represents the digit returned by the "let\_dig" module FIG. 11 that represents the nucleotide being examined.  $ALF$  is a constant that is set by the program in this implementation to equal four. The query pointer is then incremented, while the size of  $k\_tuple$  (the seed) is decremented 292. This process is repeated until the sequence of  $k\_tuple$  has been entirely hashed. Then the "tran" module returns the variable  $current\_hash$  293 to the "hashing" module FIG. 9.

The "let\_dig" module, FIG. 11, is called by the "tran" module 291, and transforms the nucleotides represented as the characters "A", "T", "U", "G" and "C" in the GenBank and the user's query into numeric digits for easier processing by the program. This module transforms "a" and "A" into "0" 301, "t", "T", "u" and "U" into "1" 302, "g" and "G" into "2" 303, and "c" and "C" into "3" 305. If the character to be transformed does not match any one of those listed above, the module returns "-1" 305. The "hashing" module, FIG. 9, then calls the "update" module 272, FIG. 12, which updates the hash with a sliding window (i.e., it forms a new hash after shifting the old hash by "1"). The remainder of  $old\_hash$  divided by  $power\_1$  is calculated 311 (a modulus operation), the remainder is multiplied by  $ALF$  312 (i.e., four), and then the digit representing the nucleotide is added to the result 313. The "update" module then returns the result 314 to the "hashing" module FIG. 9.

If the current hash has already occurred in the query, the program searches for the end of the linked list for the current hash 273 and marks the end of the linked list for the current hash 274. If the current hash has not already occurred in the query, the program puts the hash into the hash table 275. The resulting hash table and linked list are then returned to the "k\_diff" module, FIG. 8 at 258.

The "assembly" module, FIG. 13, extracts sequences from the GenBank and performs hit locating and extending functions. This module is called by the "k\_diff" module FIG. 8 at 260 if the user has chosen to use the database to locate matches. The output from the "assembly" module (FIG. 13) tells the user that the section of the database searched contains  $E$  number of entries 321 of  $S$  summary length 322 with  $H$  number of hits 323. Further, the program tells the user that the number of considered  $l$ -tuples equals  $T$  324. The entry head line is also printed 326.

26

The "seqload" module, FIG. 14, is called by the "k\_diff" module FIG. 8 at 259 once the query hash table and linked list have been formed by the "hashing" module FIG. 9. The "seqload" module FIG. 14 checks to see if the end of the GenBank file has been reached 327, and, if not, searches until a record is found with LOCUS in the head-line 328. Next, the LOCUS name is extracted 329 for identification purposes, and the program searches for the ORIGIN field in the record 330.

The program then extracts the current sequence 331 from the GenBank and performs two passes on each sequence. The first is to determine the sequence length 332 and allocate memory for each sequence 333, and the second pass is to read the sequence into the allocated memory 334. Since the sequences being extracted can contain either DNA nucleotides or protein nucleotides, the "seqload" module can recognize the characters "A", "T", "U", "G", and "C". The bases "A", "T", "G" and "C" are used in DNA sequences, while the bases "A", "U", "G" and "C" are used in RNA and mRNA sequences. The extracted sequence is then positioned according to the type of nucleotides contained in the sequence 335, and the process is repeated. Once the end of the sequence has been reached, the "seqload" module returns the sequence length 336 to the "k\_diff" module FIG. 8.

If the user has chosen to use one or more files to locate matches, rather than the database, the "read1" module, FIG. 15, rather than the "seqload" module FIG. 14, is called by the "k\_diff" module FIG. 8. The "read1" module, FIG. 15, reads the sequence from the user specified query file 341 and allocates memory 342. This module also determines the query length 343, extracts sequence identification information 344, determines the sequence length 345, transforms each nucleotide into a digit 346 by calling the "let\_dig" module FIG. 11, creates the query hash table 347 by calling the "dig\_let" module FIG. 16, and closes the file 348 once everything has been read in.

First, the "read1" module FIG. 15 allocates space for the query 342. To do this, the "ckalloc" module, FIG. 15 at 342, is called. This module allocates space and checks whether this allocation is successful (i.e., is there enough memory or has the program run out of memory). After allocating space, the "read1" module FIG. 15 opens the user-specified file 349 (the "ckopen" module, FIG. 15 at 349, is called to ensure that the query file can be successfully opened 349), determines the query length 343, locates a record with LOCUS in the head-line and extracts the LOCUS name 344 for identification purposes, locates the ORIGIN field in the record and then reads the query sequence from the file 341. Next, the sequence length is determined 345, memory is

27

allocated for the sequence 342, and the sequence is read into the query file 350. If the string has previously been found, processing is returned to 344. If not, then each character in the query file is read into memory 350.

The characters are transformed into digits 346 using the "let\_dig" module, FIG. 11, until a valid digit has been found, and then the hash table containing the query is set up 347 using the module "dig\_let", FIG. 16, which transforms the digits into nucleotides represented by the characters "A"371, "T"371, "G"373, "C"374, and "X"375 as a default. If the end of the file has not been reached, processing is returned to 344. If it has, the file is closed 348 and the query is then returned to the "read1" module FIG. 15 at 347.

The "q\_colour" module, FIG. 17 (FIG. 13 at 325), is called by the "assembly" module FIG. 13 after the current sequence has been extracted from the GenBank. The "q\_colour" module FIG. 17 performs the heart of the Mismatch Model process in that it performs the comparison between the query and the database or file sequences. If the module finds that there exists a long (i.e., greater than the min\_hit\_length) extended hit, it returns a "1" to the "assembly" module FIG. 14. Otherwise, the "q\_colour" module, FIG. 17, returns a "0".

In the "q\_colour" module, FIG. 17, all DNA positions are analyzed in the following manner. First, the entire DNA sequence is analyzed 391 to see whether each position is equal to zero 392 (i.e., whether it is empty or the sequence is finished). If it is not equal to zero 393, the "q\_colour" module FIG. 20 calls the "tran" module, FIG. 10 described above, which performs the hashing of k\_tuples. The "tran" module FIG. 10 calls other modules which transform the nucleotides represented by characters into digits for easier processing by the program and then updates the hash with a sliding window. If the position is equal to zero, the current\_hash position is set to new\_has after one shift of old\_hash 390 by calling the "update" module FIG. 12.

If the nucleotide at the current\_hash position is equal to zero, processing is returned to 391. If not, the query position is set equal to (nucleotide at current hash position - 1). Next, the "q\_colour" module FIG. 17 looks for the current\_hash in the hash table 394. If the current k\_tuple does not match the query 395, then the next k\_tuple is considered 395, and processing is returned to 391. If the current k\_tuple does match the query, then the program checks the hit's (i.e., the match's) vicinity 396 by calling the "hit\_ext" module, FIG. 18 to determine if the hit is weak. The inventors have found that if the code for the module "hit\_ext" is included within the module "q\_colour",

28

rather than being a separate module utilizing the parameter transfer machinery, 25% of CPU time can be saved.

The "hit\_ext" module FIG. 18 determines the current query position in the hit's vicinity 421, determines the current DNA position in the hit's vicinity 422, and creates the list of mismatch positions (i.e., the mismatch\_location\_ahead 423, the mismatch\_location\_behind 423 and the kernel match location). If the hit is weak 424, the "hit\_ext" module FIG. 18 returns "0" to the "q\_colour" module FIG. 17. If the hit has a chance to contain 425, the module returns "1" to the "q\_colour" module FIG. 17. A hit has a chance to contain, and is therefore not considered weak, if the mismatch\_location\_ahead - the mismatch\_location\_behind is greater than the min\_hit\_length. If not, it is a short hit and is too weak.

If the "hit\_ext" module FIG. 18 tells the "q\_colour" module FIG. 17 that the hit was not a weak one, then the "q\_colour" module determines whether the current hit is long enough 398 by calling the "colour" module FIG. 19. The "colour" module FIG. 19 performs query\_colour modification by the hit data, starting at pos\_query and described by mismatch\_location\_ahead and mismatch\_location\_behind. After the variables to be used in this module are defined, variable isw\_print (which is the switch indicating the hit length) is initialized to zero 430. The cur\_length is then set equal to the length of the extending hit 431 (mismatch\_location\_behind[i] + mismatch\_location\_ahead[j]-1). Next, if cur\_length is greater than or equal to the min\_hit\_length 432 (i.e., the minimum considered probe size), the hit is considered long and isw\_print is set equal to two 433. The value of isw\_print is then returned 434 to the "q\_colour" module FIG. 17.

If the length of the extending hit is longer than the min\_hit\_length, the hit is considered long 399. Otherwise, the hit is considered short. If the hit is short, nothing more is done to the current hit and the module begins again. If, on the other hand, the hit is considered long 399, the "q\_colour" module FIG. 17 prints the current extended hit 400. The current extended hit can be printed in ASCII, printed in a binary file, or printed to a memory file. The "q\_colour" module FIG. 17 then repeats until the end of the linked list is reached.

#### d. Outputs

The output of the k\_diff program in the current implementation of this invention may be either a binary file containing the number of extended hits and the k\_mismatch hit locations (see FIG. 20), or the output may be kept in memory without writing it to a file. See Section 1(d)(iv) for more detail.

29

### 3. Description of the H-Site Model Program

#### a. **Overview**

In this invention, the second hybridization strength model is termed the H-Site Model (see FIG. 2 for user selection of this model). One aspect of the H-Site Model uses a generalization of an experimental formula in general usage. The formula used in the H-Site Model is an expression of the fact that melting temperature  $T_m$  is a function of both probe length and percent of GC content. This basic formula has been modified in this invention to account for the presence of mismatches. Each percent of mismatch reduces the melting temperature  $T_m$  by an average of 1.25 degrees (2 degrees C for an AT mismatch, and 4 degrees C for a GC mismatch).

In addition, this implementation of the invention does some preliminary preprocessing of the GenBank database to sort out and select the cDNA sequences. This is done by locating a keyword (in this case CDS) in each GenBank record. No other programs currently available allow for this combination of functions as far as the inventors are aware.

There are a number of modules in the present embodiment of the H-Site Model contained in this invention. Each step of the processing involved in the H-Site Model is more fully explained below, and is accompanied by detailed flow charts.

#### b. **Inputs**

There are two basic user-selected inputs for the H-Site Model (see FIG. 2C): 1) the melting temperature  $T_m$  22 for which probes are being designed (i.e., the melting temperature that corresponds to a particular experiment or condition the user desires to simulate); and 2) the nucleation threshold 23, which is the number of base pairs constituting a nucleation site. The user is also required to select the 1) target species 11 gene sequence(s) (DNA, mRNA or cDNA) for which probes are being designed; 2) the preparation 12 of all sequences with which hybridizations are to be calculated; and 3) the probe output file 13. The preparation file is the most important, as discussed below.

#### c. **Organization of the H-Site Model Program**

The current implementation of the H-Site Model program of this invention is distributed between five files containing numerous modules. The main file is designated by the inventors as "ds.cpp" in its uncompiled version. This file provides overall control to the entire invention. It is divided into six sections. Section 0 defines and manipulates global variables. Section 1 controls general variable definition and initialization

(including the arrays and memory blocks). It also reads and writes buffers for user input selections, and constructs multi buffers.

Section 2 sets up and initializes various "snippet" variables (see section below for a complete definition of the term snippet), converts base pair characters to a representation that is 96 base pairs long and to ASCII base pair strings, and performs other sequence file manipulation such as comparing snippets. This section also reads the sequence format file, reads base pairs, checks for and extracts sequence identification information (such as ORIGIN and LOCUS) and filters out sequences beginning with numbers.

Section 3 involves preparation file manipulation. This section performs the preprocessing on the PRP file discussed above. It also merges and sorts the snippet files, creates a PRP file and sorts it, and outputs the sorted snippets. Next, this section streams through the PRP file.

Section 4 contains the essential code for H-Site Model processing (see FIGS. 21 through 23 for details, discussed below). Streams are set up, and then RIBI comparisons are performed for hybridizations (see file "ribi.cpp" for definitions of RIBI search techniques). Next, probes are generated, binding strength is converted to melting temperature, and hybridizations are calculated and stored (including hybridization strength). Lastly, other H-Site calculations are performed.

Section 5 is concerned with formatting and presenting diagnostic and user file (test.out, test1.out, and test2.out files) output. This section also handles the graphing functions (the MPSD diagram in particular). In addition, this section calculates the hairpin characteristics for the H-Site Model candidate probes.

The second H-Site Model file, designated as "ds.h" defines data variables and structures. Section 1 of this file concerns generic data structures (including memory blocks and arrays, and file inputs and outputs). Section 2 defines the variables and structures used with sequences, probes and hybridizations. Section 3 defines variables and structures concerned with protocols (i.e., function prototypes, graphing, etc.).

The third H-Site Model file, designated as "funcdoc.txt", contains very detailed documentation for this implementation of the H-Site Model program. Numerous variables and structures are also defined. The flow of the program is clearly shown in this file.

The fourth H-Site Model file, designated as "ribi.h" handles the sequence comparisons. The fifth and last H-Site Model file, designated as "ribi.cpp", performs



internal B-Tree indexing. Definitions of Red-black Internal Binary Index (RIBI) searching are found in this file. Definitions are also included for the concepts keyed set, index, binary tree, internal binary index, paths, and red-black trees. Implementation notes are also included in this file.

d. **Processing**

Implementation of the H-Site Model in this invention is done in three stages. First, the invention creates the preparation (PRP) file, which contains all relevant information from the sequence database. This is the preprocessing stage discussed above. Next, the target is prepared by the program. Lastly, the invention calculates the MPSD data using the PRP file and target sequence to find probes.

i. Creation of the Preprocessed Preparation File

FIG. 21. Step 1: The program first opens the sequence database for reading into memory 461, 462. Step 2: Next, as sequence base pairs are read in 462, "snippets" are saved to disk 463, along with loci information. A snippet is a fixed-length subsequence of a preparation sequence. The purpose of snippets is to allow the user to examine a small portion of a preparation sequence together with its surrounding base pairs. Snippets in the implementation of this invention are 96 base pairs long (except for snippets near the end or beginning of a sequence, which may have fewer base pairs). The "origin" of the snippet is in position 40. For snippets taken near the beginning of a sequence, some of the initial 40 bases are undefined. For snippets near the end of a sequence, some of the final 55 bases are undefined. Snippets are arranged in the preparation file (PRP) in sorted order (lexicographical order beginning at position 40). In this invention, the term "lexicographical order" means a preselected order, such as alphabetical, numeric or alphanumeric. In order to conserve space, snippets are only taken at every 4th position of the preparation sequence.

Step 3: The snippets are merge sorted 464 to be able to search quickly for sequences which pass the "screen", discussed below. Step 4: The merged file is prepended with identifiers for the sources of the snippets 465. This is done to identify the loci from which hybridizations arise.

ii. Target Preparation

FIG. 22. Step 1: The target sequence file is opened 471 and read into memory 472. For each position in the target mRNA, the probe defined at that starting position is the shortest subsequence starting at that position whose hybridization strength is greater than the user specified melting temperature  $T_m$ . Typically, the probes are of

32

length 18 to 50. Step 2: Four lists of "screens" are formed 473, 474, 475, each shifted by one base pair 475 to correspond to the fact that snippets are only taken at every four base pairs. A screen is a subsequence of the target mRNA of length equal to the screening threshold specified by the user. The screens are then indexed 476 and sorted in memory 477.

iii. Calculation of the MPSD Data

FIG. 23. Step 3: This step is the heart of the process. Step 3a: The program streams through the following five items in sync, examining them in sequential order: the snippet file and the four lists of screens 481-484. Step 3b: Each snippet is compared to a screen 485. Step 3c: If the snippet does not match, whichever stream is behind is advanced 486 and Step 3b is repeated. If the snippet does match, Step 4 is performed.

Step 4: If a snippet and a matching screen were found in Step 3b 487, the hybridization strength of the binding between the sequence containing the snippet and all of the probes containing the screen is calculated (see Step 5). Double counting is avoided by doing this only for the first matched screen containing the probe. Each pair of bases is examined and assigned a numerical binding strength. An AT pair would be assigned a lower binding strength than a GC pair because AT pairs have a lower melting temperature  $T_m$ . The process is explained more fully below at Step 5b.

Step 5: The hybridization strengths between sequence and all the probes containing it are calculated using a dynamic programming process. The process is as follows: Step 5a: Begin at the position of the first probe containing the given screen but not containing any other screens which start at an earlier position and also match the sequence. This is done to avoid double counting. Two running totals are maintained: a) boundStrength, which represents the hybridization strength contribution which would result if the sequence and probe were to match exactly for all base pairs to the right of the current position, and b) unboundStrength, which represents the strength of the maximally binding region. Step 5b: At each new base pair, the variable boundStrength is incremented by 71 if the sequence and probe match and the matched base pair is GC 489, incremented by 30 if the matched base pair is AT 490 (i.e., this number is about 42.25% of the first number 71), and decremented by 74.5 if there is not a match 488 (i.e., this number is about 5% larger than the first number 71). Step 5c: If the current boundStrength exceeds the current unboundStrength 491 (which was originally initialized to zero), a new binding region has been found, and

33

unboundStrength is set equal to boundStrength 492. Step 5d: If the current boundStrength is negative, boundStrength is reset to zero 493. Step 5e: If the current position is at the end of a probe, the results (the hybridization strengths) are tallied for that probe. Step 5f: If the current position is at the end of the last probe containing the screen, the process stops.

Step 6: A tally is kept of the number and melting temperature of the matches for each candidate probe, and the location of the best 20 candidates, using a priority queue (reverse order by hybridization strength number) 494. Step 7: A numerical "score" is kept for each preparation sequence by tallying the quantity  $\exp(-T_m)$  (which can be expressed as  $\Sigma e^{-T_m}$ ) for each match 495, where  $T_m$  is the melting temperature for the "perfect" match, the probe itself. In other words, the probe hybridizes "perfectly" to its target.

Step 8: Hairpins are calculated by first calculating the complementary probe. In other words, the order of the bases in the candidate probe are reversed (CTATAG to GATATC), and complementary base pairs are substituted (A for T, T for A, G for C, and C for G, changing GATATC to CTATAG in the above example). Next, the variable representing the maximum hairpin length for a candidate probe is initialized to zero, as is the variable representing a hairpin's distance. For each offset, the original candidate probe and the complementary probe just created are then aligned with each other and compared. The longest match is then found. If any two matches have the same length, the one with the longest hairpin distance (i.e., the number of base pairs separating the match) is then saved.

Step 9: The preparation sequences are then sorted 496 and displayed in rank order, from best to worst 497. Step 10: The resulting MPSD, which includes all candidate probes, is then displayed on the screen. Step 11: The best 20 matches are also printed or displayed in rank order, as the user requests 497.

#### e. Outputs

The outputs of the H-Site Model as currently implemented in this invention are fully described in Section 1(d)(iv), above, and illustrated in FIGS. 4 through 6. Samples of the two output files created by the H-Site Model are shown in FIGS. 24A and 24B.

#### 4. Description of the Mitsuhashi Probe Selection Diagram Processing

Once the Mitsuhashi Probe Selection Diagram (MPSD) data has been calculated by the H-Site Model program (see stage three and FIG. 23, discussed above), it is

34

necessary to convert this data to pixel format and plot a graph. An overview of this process is shown in FIG. 25. First, the program calculates the output (x,y) ranges 500. Next, these are converted to a logarithmic scale 501. The values are then interpolated 502, and a bitmap is created 503. Lastly, the bitmap is displayed on the screen 504 in MPSD format (discussed above in section 1(e)(i)). A sample MPSD is shown in FIG. 4.

#### 5. Description of the MatchInfo Window Processing

The ProbeInfo and MatchInfo windows are discussed in great detail in Section 1(e)(ii), and a sample of these windows is shown in FIG. 5. An overview of the processing involved in creating the MatchInfo portion of the window is given in the flow chart in FIG. 26. First, as the user moves the MPSD cursor 520 (seen as a vertical line bisecting the MPSD window), the program updates the position of the candidate probe shown under that cursor position 521. Next, based upon the candidate probe's position, the program updates the sequence 522 and hairpin information 523 for that probe. This updated information is then displayed in an updated match list 524, shown in the MatchInfo window.

The above described embodiments of the present invention are merely descriptive of its principles and are not to be considered limiting. The scope of the present invention instead shall be determined from the scope of the following claims including their equivalents.

35

**WHAT IS CLAIMED IS:**

1. A programmed computer system for designing optimal oligonucleotide sequences for use with a gene sequence data source comprising:

first input means for introducing user-selected gene sequence into the computer system;

memory means for storing user-selected gene sequence;

means for accessing gene sequence data from said gene sequence data source;

means for performing exact and inexact match modeling between gene sequences;

means for performing hybridization strength modeling on gene sequences;

means for selecting either of said modeling means; and

means for presenting the results of said modeling to present candidate oligonucleotide sequences.

2. A programmed computer system in accordance with Claim 1 wherein said means for performing exact and inexact match modeling utilizes said accessing means to introduce a user-selected set of gene sequence data and a user-selected set of target gene sequence data from said gene sequence data source into the computer system and said memory means to store said gene sequence data and said target gene sequence data and wherein said means for performing exact and inexact match modeling includes:

means for determining a minimum sequence length;

means for creating a look-up hash table and linked list in memory for each gene sequence in said gene sequence data and each of said target gene sequences;

means for calculating the minimum length of any matching gene subsequence of said gene sequence data and said target gene sequence data;

means for comparing each base pair character in each said target sequence stored in a hash table in memory to each base pair character of said gene sequence stored in a hash table in memory;

means for finding a matching seed by determining if the said comparison results in a matching gene subsequence of length equal to said calculated minimum length;

36

means for comparing base pair characters behind and ahead of said seed to determine if there exists an extended match of a subsequence of base pair characters of length greater than the calculated minimum length, resulting in a current hit sequence;

means for calculating whether said current hit sequence is longer than said minimum sequence length, resulting in a current candidate oligonucleotide sequence;

means for storing said current candidate oligonucleotide sequence; and

wherein said presenting means provides said current candidate oligonucleotide sequence to the user.

3. A programmed computer system in accordance with Claim 2 wherein said computer system includes:

means for calculating the melting temperature for each candidate oligonucleotide sequence;

means for tracking the number and melting temperature of the matches for-each candidate oligonucleotide sequence;

means for tracking the location of a set number of the best candidate oligonucleotide sequences; and

wherein said presenting means is operative to present said additional results to the user; and

wherein said presenting means provides said melting temperature to the user.

4. A programmed computer system in accordance with Claim 2 wherein said computer system includes:

means for determining the length of sequences from said target gene sequence data.

5. A programmed computer system in accordance with Claim 2 wherein said computer system includes:

means for determining the length of sequences from said set of gene sequence data.

37

6. A programmed computer system in accordance with Claim 2 wherein said computer system includes:

means for copying the LOCUS name for each said gene sequence into said memory means; and

means for linking said LOCUS name with each said gene sequence.

7. A programmed computer system in accordance with Claim 2 wherein said means for performing exact and inexact match modeling utilizes said accessing means to introduce a user-selected minimum sequence length from said gene sequence data source into the computer system and said memory means to store said minimum sequence length.

8. A programmed computer system in accordance with Claim 2 wherein said computer system includes:

means for calculating the melting temperature for each candidate oligonucleotide sequence;

means for tracking the number and melting temperature of the matches for each candidate oligonucleotide sequence;

means for tracking the location of a set number of the best candidate oligonucleotide sequences employing a priority queue by sorting said candidate oligonucleotide sequences in reverse order and sorting said candidate oligonucleotide sequences by hybridization strength;

wherein said presenting means is operative to present said additional results to the user; and

wherein said presenting means provides said melting temperature to the user.

9. A programmed computer system in accordance with Claim 2 wherein said first input means is operative to introduce a user-selected maximum number of mismatches and a user-selected minimum candidate oligonucleotide sequence length into the computer system, and wherein said means for calculating the minimum length of any matching gene subsequence of said gene sequence data and said target gene sequence data comprises the steps of:

38

means for subtracting said maximum number of mismatches from said minimum candidate oligonucleotide sequence length to give a first result;

means for dividing said first result by said maximum number of mismatches plus one to give a second result;

means for incrementing said second result by one if the remainder is not equal to zero to give a third result; and

means for truncating said third result to an integer.

10. A programmed computer system in accordance with Claim 9 wherein said means for calculating the hairpin characteristics of said candidate oligonucleotide sequence comprises the steps of:

calculating a complementary sequence to the candidate oligonucleotide sequence by reversing the base pair order of the candidate oligonucleotide sequence and substituting complementary base pairs;

comparing each character of said original candidate oligonucleotide sequence and said complementary sequence;

finding the longest match between said original candidate oligonucleotide sequence and said complementary sequence; and

saving the match with the longest hairpin distance if any two matches have the same length;

means for storing hairpin characteristics; and

wherein said presenting means provides said hairpin characteristics to the user.

11. A programmed computer system in accordance with Claim 2 wherein said computer system includes a means for calculating the hairpin characteristics of said candidate oligonucleotide sequence.

12. A programmed computer system in accordance with Claim 2 wherein said means for preprocessing said set of target gene sequence data and said set of gene sequence data comprises the steps of:

searching for sequences without introns in said target gene sequences and said gene sequences;



39

extracting target gene sequences and gene sequences that do not contain introns; and

storing said extracted target gene sequences and gene sequences in memory.

13. A programmed computer system in accordance with Claim 1 wherein said means for performing hybridization strength modeling utilizes said first input means to introduce a user-selected screening threshold into the computer system and said accessing means to introduce a user-selected set of gene sequence data and a user-selected set of target gene sequence data from said gene sequence data source into the computer system, and said memory means to store said gene sequence data, said target gene sequence data and said screening threshold and wherein said means for performing hybridization strength modeling comprises:

means for preprocessing said target gene sequence data and said gene sequence data by selecting only those sequences without introns;

means for forming a preparation file of gene sequence fragments by cutting said target gene sequences into fixed length target gene subsequences and sorting said subsequences in lexicographical order;

means for merge sorting said gene sequences;

means for forming multiple lists of screens by forming lists of subsequences of the preparation file of length equal to said screening threshold;

means for indexing, sorting and storing said screens in said memory means;

means for sequentially comparing said preparation file gene sequences with each of said screens to design candidate oligonucleotide sequences;

means for calculating the hybridization strengths between a gene sequence and all candidate oligonucleotide sequences containing that gene sequence by accounting for Guanine-Cytosine (GC) and Adenine-Thymine (AT) base pair content of the gene sequence and the number of mismatches between said preparation file sequences and a said screen when said comparison results in a match;

means for preparing the candidate oligonucleotide sequence and hybridization strength for presentation to the user; and

wherein said presenting means provides the candidate oligonucleotide sequence and hybridization strength to the user.

40

14. A programmed computer system in accordance with Claim 13 wherein said computer system includes:

means for calculating the melting temperature for each candidate oligonucleotide sequence;

means for tracking the number and melting temperature of the matches for each candidate oligonucleotide sequence;

means for tracking the location of a set number of the best candidate oligonucleotide sequences;

means for preparing the melting temperature for presentation to the user; and

wherein said presenting means provides the melting temperature to the user.

15. A programmed computer system in accordance with Claim 14 wherein said means for calculating said candidate oligonucleotide sequence's melting temperature comprises:

solving the formula  $T_m = 81.5 - 16.6(\log[Na]) - .63 \%(formamide) + ((.41 (\%(G + C)) - 600)/N)$ , wherein  $\log[Na]$  is the sodium concentration,  $\%(G + C)$  is the fraction of matched base pairs which are G-C complementary, N is the sequence length and wherein the number of mismatches is equal to zero.

16. A programmed computer system in accordance with Claim 15 wherein said computer system includes:

means for reducing a candidate oligonucleotide probe's calculated melting temperature by a certain amount for each percent of mismatch between the candidate oligonucleotide sequence and a user-selected target gene sequence based upon the assumption that there are an equal number of GC and AT base pair mismatches.

17. A programmed computer system in accordance with Claim 16 wherein said means for reducing a candidate oligonucleotide sequence's calculated melting temperature comprises the steps of:

reducing said calculated melting temperature by 2 degrees Celsius if an AT mismatch exists; and

41

reducing said calculated melting temperature by 4 degrees Celsius if a GC mismatch exists.

18. A programmed computer system in accordance with Claim 13 wherein said computer system includes:

means for assigning a numerical score to each said gene sequence; and  
means for sorting said gene sequences in accordance with said numerical score.

19. A programmed computer system in accordance with Claim 13 wherein said means for performing hybridization strength modeling utilizes said accessing means for copying the LOCUS name for each said gene sequence into said memory means, and said memory means; and

means for prepending said gene sequence with said LOCUS name.

20. A programmed computer system in accordance with Claim 13 wherein four lists of screens are formed by said list forming means.

21. A programmed computer system in accordance with Claim 13 wherein said computer system includes a means of shifting each screen by at least one base pair as it is formed by said list forming means.

22. A programmed computer system in accordance with Claim 13 wherein said computer system includes:

means for calculating the melting temperature for each candidate oligonucleotide sequence;

means for tracking the number and melting temperature of the matches for each candidate oligonucleotide sequence;

means for tracking the location of a set number of the best candidate oligonucleotide sequences employing a priority queue by sorting said candidate oligonucleotide sequences in reverse order and sorting said candidate oligonucleotide sequences by hybridization strength;

means for preparing the melting temperature for presentation to the user;  
and

42

wherein said presenting means provides the melting temperature to the user.

23. A programmed computer system in accordance with Claim 13 wherein said computer system includes:

means for assigning a numerical score to each said gene sequence by tallying the quantity "exp" where  $\text{"exp"} = \Sigma e^{-T_m}$  and wherein  $T_m$  is the melting temperature for the said gene sequence; and

means for sorting said gene sequences in accordance with said numerical score.

24. A programmed computer system in accordance with Claim 13 wherein said means for calculating the hybridization strengths between a gene sequence and all candidate oligonucleotide sequences containing that gene sequence comprises the steps of:

accessing gene sequence data from said gene sequence data source;

comparing base pairs of a first gene sequence and a second gene sequence to determine if a match exists;

incrementing said first gene sequence's bound strength by some first number if a base pair character in said first gene sequence and said second gene sequence match and the matched base pair is equal to a combination of the bases Guanine (G) and Cytosine (C);

incrementing said first gene sequence's bound strength by some second number if a base pair character in said first gene sequence and said second gene sequence match and the matched base pair is equal to a combination of the bases Adenine (A) and Thymine (T);

decrementing said first gene sequence's bound strength by a third number if there is no match in base pairs between said first gene sequence and said second gene sequence;

comparing said first gene sequence's bound strength to said first gene sequence's unbound strength;

setting said first gene sequence's unbound strength equal to its bound strength if said first gene sequence's bound strength is greater than said first gene sequence's unbound strength; and

43

resetting said first gene sequence's bound strength to zero if said first gene sequence's unbound strength is less than zero.

25. A programmed computer system in accordance with Claim 24 wherein said first and second numbers are greater than zero.

26. A programmed computer system in accordance with Claim 24 wherein said second number is in the order of 42% of said first number.

27. A programmed computer system in accordance with Claim 24 wherein said third number is in the order of 5% larger than said first number.

28. A programmed computer system in accordance with Claim 13 wherein said computer system includes a means for calculating the hairpin characteristics of said candidate oligonucleotide sequence;

means for preparing the hairpin characteristics for presentation to the user;

and

wherein said presenting means provides the hairpin characteristics to the user.

29. A programmed computer system in accordance with Claim 28 wherein said means for calculating the hairpin characteristics of said candidate oligonucleotide sequence comprises the steps of:

calculating a complementary sequence to the candidate oligonucleotide sequence by reversing the base pair order of the candidate oligonucleotide sequence and substituting complementary base pairs;

comparing each character of said original candidate oligonucleotide sequence and said complementary sequence;

finding the longest match between said original candidate oligonucleotide sequence and said complementary sequence; and

saving the match with the longest hairpin distance if any two matches have the same length;

means for preparing the hairpin characteristics for presentation to the user;

and

44

wherein said presenting means provides the hairpin characteristics to the user.

30. A programmed computer system in accordance with Claim 13 wherein said fixed-length subsequences are calculated by a method comprising the steps of:

locating the origin of said subsequence in a set position of said target gene sequence in said preparation file;

cutting a subsequence that is a fixed-length long every preselected number of positions of said target gene sequence in said preparation file; and

sorting said subsequences in said preparation file in lexicographical order beginning at a set position.

31. A programmed computer system in accordance with Claim 30 wherein the origin of said subsequence is located at position 40 of said target sequence in said preparation file.

32. A programmed computer system in accordance with Claim 13 wherein said fixed-length subsequences are calculated by a method comprising the steps of:

locating the origin of said subsequence in the 40th position of said target gene sequence in said preparation file;

cutting a subsequence that is 96 base pairs long of said target gene sequence in said preparation file; and

sorting said subsequences in said preparation file in lexicographical order beginning at a set position.

33. A programmed computer system in accordance with Claim 13 wherein said computer system includes means for prepending said preparation file subsequences with identifiers for the sources of each subsequence.

34. A programmed computer system in accordance with Claim 1 wherein said presenting means to provide the results of said matching and modeling to display candidate oligonucleotide sequences includes means for displaying in multiple dimensions the gene sequences which result from the comparisons and calculations characterized in that said display format exhibits

45

the starting position of each candidate oligonucleotide sequence in one dimension;

the specificity of a candidate oligonucleotide sequence's hybridization with the target gene sequence in a second dimension; and

superimposed melting temperatures of gene sequences in contrasting presentations in at least an apparent third dimension.

35. A programmed computer system in accordance with Claim 34 wherein said display further includes a cursor moveable along one dimension of said display that selects a position for an expansion of data representing the homology between the candidate oligonucleotide sequences and said gene sequence data; and

wherein said display is operative to display in alphanumeric form the homology between the candidate oligonucleotide sequences and said gene sequence data.

36. A programmed computer system in accordance with Claim 34 wherein said display is further operative to provide an expansion of data including presenting

false hybridizations at various melting temperatures

for all candidate oligonucleotide sequences;

the location of each false hybridization;

a candidate oligonucleotide sequence's starting position; and

hairpin characteristics of each candidate oligonucleotide sequence.

37. A programmed computer system in accordance with Claim 34 wherein said display format data is outputted to a printing means.

38. A programmed computer system in accordance with Claim 34 wherein said display format data is saved to a data file.

39. A programmed computer system in accordance with Claim 34 wherein said display format data is exported to another computer system.

4/6

40. A programmed computer system in accordance with Claim 34 wherein said display further includes a cursor moveable along one dimension of said display that selects a position for an expansion of data representing the homology between the candidate oligonucleotide sequences and said gene sequence data; and

wherein said moveable cursor may be positioned by the user to select and save particular candidate oligonucleotide sequence information; and

wherein said display is operative to display in alphanumeric form the homology between the candidate oligonucleotide sequences and said gene sequence data.

41. A programmed computer system in accordance with Claim 40 wherein said method of selecting and saving particular candidate oligonucleotide sequence information comprises capturing candidate oligonucleotide sequence information at the user-selected point and storing said information in said memory means.

42. A programmed computer system in accordance with Claim 41 wherein said user-selected candidate oligonucleotide sequence information is exported to another computer system.

43. A programmed computer system in accordance with Claim 34 wherein said means for displaying comprises the steps of:

- calculating display output ranges;
- converting said output ranges to a logarithmic scale;
- interpolating said converted values;
- creating a bitmap of said interpolations; and
- displaying said bitmap on a display device.

44. A programmed computer system in accordance with Claim 34 wherein said means for displaying comprises the steps of:

- converting said result values to pixels;
- filling a pixel array with said pixels;
- performing a binary search into said pixel array;
- determining the number of pixels per candidate oligonucleotide sequence to be displayed;



47

interpolating said pixels at the value of pixels per position minus one;  
computing an array of said pixel array; and  
plotting the results on a display device.

45. A programmed computer system in accordance with Claim 1 wherein said means for performing exact and inexact match modeling utilizes said accessing means to introduce a user-selected set of gene sequence data and a user-selected set of target gene sequence data from said gene sequence data source into the computer system and said memory means to store said gene sequence data and said target gene sequence data and wherein said means for performing exact and inexact match modeling includes:
- means for determining a minimum sequence length;
  - means for creating a look-up hash table and linked list in memory for each gene sequence in said gene sequence data and each of said target gene sequences;
  - means for calculating the minimum length of any matching gene subsequence of said gene sequence data and said target gene sequence data;
  - means for transforming base characters in each said target sequence and in each said gene sequence into numeric digits;
  - means for comparing each base pair digit in each said target sequence stored in a hash table in memory to each base pair digit of said gene sequence stored in a hash table in memory;
  - means for finding a matching seed by determining if the said comparison results in a matching gene subsequence of length equal to said calculated minimum length;
  - means for comparing base pair digits behind and ahead of said seed to determine if there exists an extended match of a subsequence of base pair digits of length greater than the calculated minimum length, resulting in a current hit sequence;
  - means for calculating whether said current hit sequence is longer than said minimum sequence length, resulting in a current candidate oligonucleotide sequence;
  - means for storing said current candidate oligonucleotide sequence; and
  - wherein said presenting means provides said current candidate oligonucleotide sequence to the user.

46. A programmed computer system for designing candidate oligonucleotide sequences for use with a gene sequence data source including:

48

first input means for introducing user-selected gene sequence, design, model and presentation criteria and a user-specified sequence length into the computer system;

memory means for storing said gene sequence, design, model and presentation criteria and said sequence length;

means for accessing gene sequence data from said gene sequence data source;

wherein said accessing means is operative to introduce a user-selected set of gene sequence data and a user-selected set of target gene sequence data from said gene sequence data source into the computer system;

wherein said criteria are used for comparison of gene sequence data and target gene sequence data;

means for comparing said gene sequences against said target gene sequences employing said criteria;

means for calculating candidate oligonucleotide sequences of said sequence-length that are either common to a pool of user-specified gene sequences or specific to a particular user-specified gene sequence;

means for calculating the homology between the candidate oligonucleotide sequences and said gene sequence data;

means for calculating a candidate oligonucleotide sequence's hairpin characteristics;

means for displaying in multiple dimensions the gene sequences which result from the comparisons and calculations characterized in that said display format exhibits:

the starting position of each candidate oligonucleotide sequence in one dimension;

a candidate oligonucleotide sequence's specificity to the target gene sequence in a second dimension; and

superimposed melting temperatures of gene sequences in contrasting presentations in at least an apparent third dimension;

wherein said display further includes a cursor moveable along one dimension of said display that selects a position for an expansion of data representing

49

the homology between the candidate oligonucleotide sequences and said gene sequence data;

wherein said display is operative to display in alphanumeric form the homology between the candidate oligonucleotide sequences and said gene sequence data; and

wherein said display is operative to provide an expansion of data including presenting

- false hybridizations at various melting temperatures
- for all candidate oligonucleotide sequences;
- the location of each false hybridization;
- a candidate oligonucleotide sequence's starting position; and
- hairpin characteristics of each candidate oligonucleotide sequence.

47. A method for designing candidate oligonucleotide sequences by performing exact and inexact match modeling for use with a gene sequence data source comprising the steps of:

- introducing user-selected gene sequence into a computer system;
- accessing gene sequence data from said gene sequence data source;
- storing user-selected gene sequence in the memory of the computer system;

- accessing the gene sequence source to introduce the user-selected set of gene sequence data and a user-selected set of target gene sequence data from said gene sequence data source into the computer system;

- storing said gene sequence data and said target gene sequence data in the memory of the computer system;

- determining a minimum sequence length;

- creating a look-up hash table and linked list in memory for each gene sequence in said gene sequence data and each of said target gene sequences;

- calculating the minimum length of any matching gene subsequence of said gene sequence data and said target gene sequence data;

50

comparing each base pair character in each said target sequence stored in a hash table in memory to each base pair character of said gene sequence stored in a hash table in memory;

determining a matching seed by determining if the said comparison results in a matching gene subsequence of length equal to said calculated minimum length;

comparing base pair characters behind and ahead of said seed to determine if there exists an extended match of a subsequence of base pair characters of length greater than the calculated minimum length, resulting in a current hit sequence;

calculating whether said current hit sequence is longer than said minimum sequence length, resulting in a current candidate oligonucleotide sequence;

storing said current candidate oligonucleotide sequence in the memory of the computer system; and

presenting a representation of said current candidate oligonucleotide sequence to the user.

48. A method in accordance with Claim 47 wherein said method includes the steps for performing additional calculations for each candidate oligonucleotide probe, said additional calculations comprising:

calculating the melting temperature for each candidate oligonucleotide sequence;

tracking the number and melting temperature of the matches for each candidate oligonucleotide sequence;

tracking the location of a set number of the best candidate oligonucleotide sequences; and

presenting said additional results to the user.

49. A method in accordance with Claim 47 wherein said method includes the step of transforming base characters into numeric digits.

50. A method in accordance with Claim 47 wherein said method includes the step of determining the length of sequences from said target gene sequence data.

51. A method in accordance with Claim 47 wherein said method includes the step of determining the length of sequences from said set of gene sequence data.

51

52. A method in accordance with Claim 47 wherein said method includes the steps of:

copying the LOCUS name for each said gene sequence into the memory of the computer system; and

linking said LOCUS name with each said gene sequence.

53. A method in accordance with Claim 47 wherein said method includes the steps of:

introducing a user-selected minimum sequence length into the computer system; and

storing said minimum sequence length in the memory of the computer system.

54. A method in accordance with Claim 47 wherein said method includes the steps for performing additional calculations for each candidate oligonucleotide probe, said additional calculations comprising:

calculating the melting temperature for each candidate oligonucleotide sequence;

tracking the number and melting temperature of the matches for each candidate oligonucleotide sequence;

tracking the location of a set number of the best candidate oligonucleotide sequences employing a priority queue by sorting said candidate oligonucleotide sequences in reverse order and sorting said candidate oligonucleotide sequences by hybridization strength; and

presenting said additional results to the user.

55. A method in accordance with Claim 47 wherein said step for calculating the minimum length of any matching gene subsequence comprises:

introducing a user-selected maximum number of mismatches and a user-selected minimum candidate oligonucleotide sequence length into the computer system;

subtracting said maximum number of mismatches from said minimum candidate oligonucleotide sequence length to give a first result;

52

dividing said first result by said maximum number of mismatches plus one to give a second result;

incrementing said second result by one if the remainder is not equal to zero to give a third result; and

truncating said third result to an integer.

56. A method in accordance with Claim 47 wherein said method includes the step of calculating the hairpin characteristics of said candidate oligonucleotide sequence.

57. A method in accordance with Claim 47 wherein said method includes the step of calculating the hairpin characteristics of said candidate oligonucleotide sequence comprising:

calculating a complementary sequence to the candidate oligonucleotide sequence by reversing the base pair order of the candidate oligonucleotide sequence and substituting complementary base pairs;

comparing each character of said original candidate oligonucleotide sequence and said complementary sequence;

finding the longest match between said original candidate oligonucleotide sequence and said complementary sequence; and

saving the match with the longest hairpin distance if any two matches have the same length.

58. A method for designing candidate oligonucleotide sequences by performing hybridization strength modeling for use with a gene sequence data source comprising the steps of:

introducing user-selected gene sequence and a user-selected screening threshold into a computer system;

storing user-selected gene sequence and said screening threshold in the memory of the computer system;

accessing the gene sequence source to introduce the user-selected set of gene sequence data and a user-selected set of target gene sequence data from said gene sequence data source into the computer system;

storing said gene sequence data and said target gene sequence data in the memory of the computer system;

53

preprocessing said target gene sequence data and said gene sequence data by selecting only those sequences without introns;

forming a preparation file of gene sequence fragments by cutting said target gene sequences into fixed length target gene subsequences and sorting said subsequences in lexicographical order;

merge sorting said gene sequences;

forming multiple lists of screens by forming lists of subsequences of the preparation file of length equal to said screening threshold;

indexing and sorting said screens in memory;

storing said screens in the memory of the computer system;

sequentially comparing said preparation file gene sequences with each of said screens to design candidate oligonucleotide sequences;

calculating the hybridization strengths between a gene sequence and all candidate oligonucleotide sequences containing that gene sequence by accounting for Guanine-Cytosine (GC) and Adenine-Thymine (AT) base pair content of the gene sequence and the number of mismatches between said preparation file sequences and a said screen when said comparison results in a match;

preparing the candidate oligonucleotide sequence and hybridization strength for presentation to the user; and

presenting the candidate oligonucleotide sequence and hybridization strength to the user.

59. A method in accordance with Claim 58 wherein said method includes the steps for performing additional calculations for each candidate oligonucleotide probe, said additional calculations comprising:

calculating the melting temperature for each candidate oligonucleotide sequence;

tracking the number and melting temperature of the matches for each candidate oligonucleotide sequence;

tracking the location of a set number of the best candidate oligonucleotide sequences; and

presenting said additional results to the user.

54

60. A method in accordance with Claim 58 wherein the step for preparing the candidate oligonucleotide sequence for presenting to the user comprises:

assigning a numerical score to each said gene sequence;  
sorting said gene sequences in accordance with said numerical score; and  
displaying a representation of the resulting candidate oligonucleotide sequence and said gene sequences.

61. A method in accordance with Claim 58 wherein said method includes the steps of:

copying the LOCUS name for each said gene sequence into the memory of the computer system; and  
prepending said gene sequence with said LOCUS name.

62. A method in accordance with Claim 58 wherein the step for forming lists of screens produces four lists of screens.

63. A method in accordance with Claim 58 wherein said method includes a the step of shifting each screen by one base pair as it is formed.

64. A method in accordance with Claim 58 wherein said method includes the steps for performing additional calculations for each candidate oligonucleotide probe, said additional calculations comprising:

calculating the melting temperature for each candidate oligonucleotide sequence;

tracking the number and melting temperature of the matches for each candidate oligonucleotide sequence;

tracking the location of a set number of the best candidate oligonucleotide sequences employing a priority queue by sorting said candidate oligonucleotide sequences in reverse order and sorting said candidate oligonucleotide sequences by hybridization strength; and

presenting said additional results to the user.

65. A method in accordance with Claim 58 wherein said method for preparing the results for presenting to the user comprises:



55

assigning a numerical score to each said gene sequence by tallying the quantity "exp" where  $\text{"exp"} = \Sigma e^{-T_m}$  and wherein  $T_m$  is the melting temperature for the said gene sequence;

sorting said gene sequences in order of the numerical score; and

displaying a representation of the resulting candidate oligonucleotide sequence and said gene sequences.

66. A method in accordance with Claim 58 for use with a gene sequence data source, programmed to determine hybridization strength comprising the steps of:

comparing base pairs of a first gene sequence and a second gene sequence to determine if a match exists;

incrementing said first gene sequence's bound strength by some first number if a base pair character in said first gene sequence and said second gene sequence match and the matched base pair is equal to a combination of the bases Guanine (G) and Cytosine (C);

incrementing said first gene sequence's bound strength by some second number if a base pair character in said first gene sequence and said second gene sequence match and the matched base pair is equal to a combination of the bases Adenine (A) and Thymine (T);

decrementing said first gene sequence's bound strength by a third number if there is no match in base pairs between said first gene sequence and said second gene sequence;

comparing said first gene sequence's bound strength to said first gene sequence's unbound strength;

setting said first gene sequence's unbound strength equal to its bound strength if said first gene sequence's bound strength is greater than said first gene sequence's unbound strength; and

resetting said first gene sequence's bound strength to zero if said first gene sequence's unbound strength is less than zero.

67. A method in accordance with Claim 66 wherein said first and second numbers are greater than zero.

56

68. A method in accordance with Claim 66 wherein said second number is in the order of 42% of said first number.

69. A method in accordance with Claim 66 wherein said second number is in the order of 5% larger than said first number.

70. A method in accordance with Claim 58 wherein said method includes the step of calculating the hairpin characteristics of said candidate oligonucleotide sequence.

71. A method in accordance with Claim 70 wherein the step of calculating the hairpin characteristics of said candidate oligonucleotide sequence includes the steps of:

calculating a complementary sequence to the candidate oligonucleotide sequence by reversing the base pair order of the candidate oligonucleotide sequence and substituting complementary base pairs;

comparing each character of said original candidate oligonucleotide sequence and said complementary sequence;

finding the longest match between said original candidate oligonucleotide sequence and said complementary sequence; and

saving the match with the longest hairpin distance if any two matches have the same length.

72. A method in accordance with Claim 58 wherein said fixed-length target gene subsequences are calculated by a method comprising the steps of:

locating the origin of said subsequence in a set position of said target gene sequence in said preparation file;

cutting a subsequence that is a fixed-length long every preselected number of positions of said target gene sequence in said preparation file; and

sorting said subsequences in said preparation file in lexicographical order beginning at a set position.

73. A method in accordance with Claim 72 wherein the origin of said subsequence is located at position 40 of said target sequence in said preparation file.

57

74. A method in accordance with Claim 58 wherein said fixed-length subsequences are calculated by a method comprising the steps of:

locating the origin of said subsequence in the 40th position of said target gene sequence in said preparation file;

cutting a subsequence that is 96 base pairs long of said target gene sequence in said preparation file; and

sorting said subsequences in said preparation file in lexicographical order beginning at a set position.

75. A method in accordance with Claim 58 wherein said method includes the step of prepending said preparation file subsequences with identifiers for the sources of each subsequence.

76. A method in accordance with Claim 58 wherein said method includes the step of calculating an candidate oligonucleotide sequence's melting temperature comprising:

solving the formula  $T_m = 81.5 - 16.6(\log[Na]) - .63 \%(formamide) + ((.41 (\%(G + C)) - 600)/N)$ ;

wherein  $\log[Na]$  is the sodium concentration,  $\%(G + C)$  is the fraction of matched base pairs which are G-C complementary, N is the sequence length; and

wherein the number of mismatches is equal to zero.

77. A method in accordance with Claim 58 wherein said method includes the step for reducing a candidate oligonucleotide sequence's calculated melting temperature by a preselected amount for each percent of mismatch between the candidate oligonucleotide sequence and a user-selected target gene sequence based upon the assumption that there are an equal number of GC and AT base pair mismatches.

78. A method in accordance with Claim 58 wherein said method includes the step for reducing a candidate oligonucleotide sequence's calculated melting temperature by a preselected amount comprising the steps of:

reducing said calculated melting temperature by 2 degrees Celsius if an AT mismatch exists; and

58

reducing said calculated melting temperature by 4 degrees Celsius if a GC mismatch exists.

79. A method for designing candidate oligonucleotide sequences for use with a gene sequence data source comprising the steps of:

introducing user-selected gene sequence and a user-specified sequence length into a computer system;

storing said gene sequence and said sequence length in the memory of the computer system;

accessing gene sequence data from said gene sequence data source;

accessing the gene sequence source to introduce the user-selected set of gene sequence data and a user-selected set of target gene sequence data from said gene sequence data source into the computer system;

comparing said gene sequences against said target gene sequences employing said criteria;

calculating candidate oligonucleotide sequences of said sequence length that are either common to a pool of user-specified gene sequences or specific to a particular user-specified gene sequence;

calculating the homology between the candidate oligonucleotide sequences and said gene sequence data;

displaying in multiple dimensions the gene sequences which result from the comparisons and calculations characterized in that said display format exhibits:

the starting position of each candidate oligonucleotide sequence in one dimension;

a candidate oligonucleotide sequence's specificity to the target gene sequence in a second dimension; and

superimposed melting temperatures of gene sequences in contrasting presentations in at least an apparent third dimension.

80. A method in accordance with Claim 79 wherein said method includes the step of calculating a candidate oligonucleotide sequence's hairpin characteristics.

81. A method in accordance with Claim 80 wherein said step of calculating hairpin characteristics for a gene sequence comprises:

59

calculating a complementary sequence to the said gene sequence by reversing the base pair order of the gene sequence and substituting complementary base pairs;

comparing each character of said original gene sequence and said complementary sequence;

finding the longest match between said original gene sequence and said complementary sequence; and

saving the match with the longest hairpin distance if any two matches have the same length.

82. A method in accordance with Claim 79 wherein the step of displaying further includes producing a cursor moveable along one dimension of said display that selects a position for an expansion of data representing the homology between the candidate oligonucleotide sequences and said gene sequence data; and

displaying in alphanumeric form the homology between the candidate oligonucleotide sequences and said gene sequence data.

83. A method in accordance with Claim 79 wherein said display format data is outputted to a printing means.

84. A method in accordance with Claim 79 wherein said display format data is saved to a data file.

85. A method in accordance with Claim 79 wherein said display format data is exported to another computer system.

86. A method in accordance with Claim 79 wherein the step of displaying further includes producing a cursor moveable along one dimension of said display that selects a position for an expansion of data representing the homology between the candidate oligonucleotide sequences and said gene sequence data;

positioning said moveable cursor to select and save particular candidate oligonucleotide sequence information; and

displaying in alphanumeric form the homology between the candidate oligonucleotide sequences and said gene sequence data.

60

87. A method in accordance with Claim 79 wherein the step of displaying further includes producing a cursor moveable along one dimension of said display that selects a position for an expansion of data representing the homology between the candidate oligonucleotide sequences and said gene sequence data;

positioning said moveable cursor to select and save particular candidate oligonucleotide sequence information;

capturing candidate oligonucleotide sequence information at the user-selected point and storing said information in said memory means; and

displaying in alphanumeric form the homology between the candidate oligonucleotide sequences and said gene sequence data.

88. A method in accordance with Claim 79 wherein said method of displaying comprises:

calculating display output ranges;

converting said output ranges to a logarithmic scale;

interpolating said converted values;

creating a bitmap of said interpolations; and

displaying said bitmap on a display device.

89. A method in accordance with Claim 79 wherein said method of displaying comprises:

converting said result values to pixels;

filling a pixel array with said pixels;

performing a binary search into said pixel array;

determining the number of pixels per candidate oligonucleotide sequence to be displayed;

interpolating said pixels at the value of pixels per position minus one;

computing an array of said pixel array; and

plotting the results on a display device.

90. A method to determine hybridization strength between two or more gene sequences for use with a gene sequence data source, comprising the steps of:

accessing gene sequence data from said gene sequence data source;

61

comparing base pairs of a first gene sequence and a second gene sequence to determine if a match exists;

incrementing said first gene sequence's bound strength by some first number if a base pair character in said first gene sequence and said second gene sequence match and the matched base pair is equal to a combination of the bases Guanine (G) and Cytosine (C);

incrementing said first gene sequence's bound strength by some second number if a base pair character in said first gene sequence and said second gene sequence match and the matched base pair is equal to a combination of the bases Adenine (A) and Thymine (T);

decrementing said first gene sequence's bound strength by a third number if there is no match in base pairs between said first gene sequence and said second gene sequence;

comparing said first gene sequence's bound strength to said first gene sequence's unbound strength;

setting said first gene sequence's unbound strength equal to its bound strength if said first gene sequence's bound strength is greater than said first gene sequence's unbound strength; and

resetting said first gene sequence's bound strength to zero if said first gene sequence's unbound strength is less than zero.

91. A method in accordance with Claim 90 wherein said first and second numbers are greater than zero.

92. A method in accordance with Claim 90 wherein said second number is in the order of 42% of said first number.

93. A method in accordance with Claim 90 wherein said third number is in the order of 5% larger than said first number.

94. A method of calculating the minimum length of any matching gene subsequence comprising:

introducing a user-selected maximum number of mismatches and a user-selected minimum candidate oligonucleotide sequence length;

62

subtracting said maximum number of mismatches from said minimum candidate oligonucleotide sequence length to give a first result;

dividing said first result by said maximum number of mismatches plus one to give a second result;

incrementing said second result by one if the remainder is not equal to zero to give a third result; and

truncating said third result to an integer.

95. A method of calculating hairpin characteristics for a gene sequence comprising:

calculating a complementary sequence to the said gene sequence by reversing the base pair order of the gene sequence and substituting complementary base pairs;

comparing each character of said original gene sequence and said complementary sequence;

finding the longest match between said original gene sequence and said complementary sequence; and

saving the match with the longest hairpin distance if any two matches have the same length.

96. A method of creating a preparation file from a user-selected set of target gene sequence data comprising:

cutting said target gene sequence data into fixed-length subsequences; and  
storing said subsequences in a preparation file.

97. A method of creating a preparation file from a user-selected set of target gene sequence data comprising:

cutting said target gene sequence data into fixed-length subsequences in the order of 96 base pairs in length; and  
storing said subsequences in a preparation file.

98. A method in accordance with Claim 97 wherein said fixed-length subsequences are calculated by a method comprising the steps of:



63

locating the origin of said subsequence in a set position of said target gene sequence in said preparation file;

cutting a subsequence that is a fixed-length long every preselected number of positions of said target gene sequence in said preparation file; and

sorting said subsequences in said preparation file in lexicographical order beginning at a set position.

99. A method in accordance with Claim 97 wherein said fixed-length subsequences are calculated by a method comprising the steps of:

locating the origin of said subsequence in a set position of said target gene sequence in said preparation file wherein the origin of said subsequence is located at position 40 of said target sequence in said preparation file;

cutting a subsequence that is a fixed-length long every preselected number of positions of said target gene sequence in said preparation file; and

sorting said subsequences in said preparation file in lexicographical order beginning at a set position.

100. A method in accordance with Claim 97 wherein said fixed-length subsequences are calculated by a method comprising the steps of:

locating the origin of said subsequence the 40th position of said target gene sequence in said preparation file;

cutting a subsequence that is 96 base pairs long of said target gene sequence in said preparation file; and

sorting said subsequences in said preparation file in lexicographical order beginning at a set position.

101. A method of forming lists of screens of target gene sequence data comprising:

introducing a user-selected screening threshold; and

forming subsequences of said target gene sequence data of length equal to a user-selected screening threshold.

102. A method of preprocessing a user-selected set of target gene sequence data comprising the steps of:

64

searching for sequences without introns in said target gene sequences;  
extracting target gene sequences that do not contain introns; and  
storing said extracted target gene sequences.

## AMENDED CLAIMS

[received by the International Bureau on 4 April 1994 (04.04.94);  
original claim 69 amended; remaining claims unchanged (1 page)]

68. A method in accordance with Claim 66 wherein said second number is in the order of 42% of said first number.

69. A method in accordance with Claim 66 wherein said third number is in the order of 5% larger than said first number.

70. A method in accordance with Claim 58 wherein said method includes the step of calculating the hairpin characteristics of said candidate oligonucleotide sequence.

71. A method in accordance with Claim 70 wherein the step of calculating the hairpin characteristics of said candidate oligonucleotide sequence includes the steps of:  
calculating a complementary sequence to the candidate oligonucleotide sequence by reversing the base pair order of the candidate oligonucleotide sequence and substituting complementary base pairs;

comparing each character of said original candidate oligonucleotide sequence and said complementary sequence;

finding the longest match between said original candidate oligonucleotide sequence and said complementary sequence; and

saving the match with the longest hairpin distance if any two matches have the same length.

72. A method in accordance with Claim 58 wherein said fixed-length target gene subsequences are calculated by a method comprising the steps of:

locating the origin of said subsequence in a set position of said target gene sequence in said preparation file;

cutting a subsequence that is a fixed-length long every preselected number of positions of said target gene sequence in said preparation file; and

sorting said subsequences in said preparation file in lexicographical order beginning at a set position.

73. A method in accordance with Claim 72 wherein the origin of said subsequence is located at position 40 of said target sequence in said preparation file.

FIG. 1

1/156

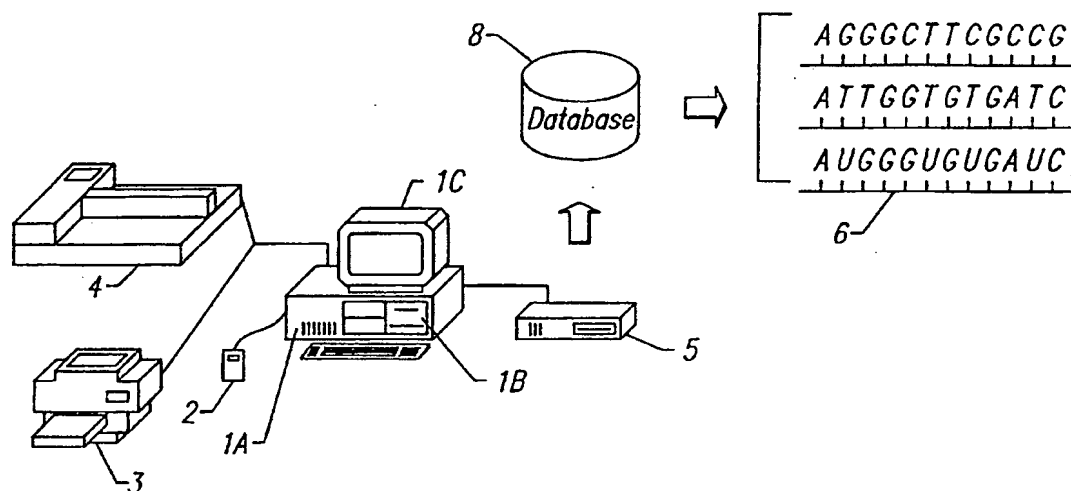
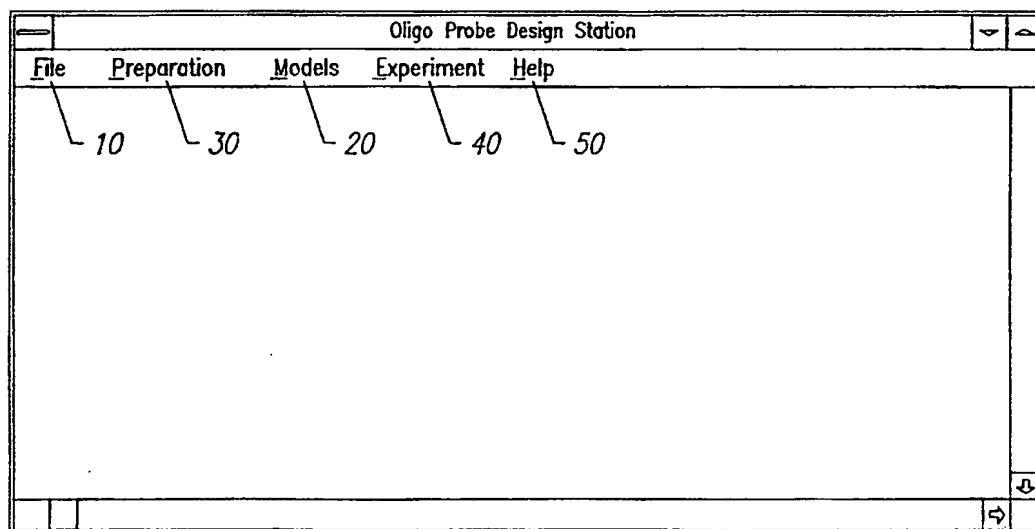


FIG. 2A



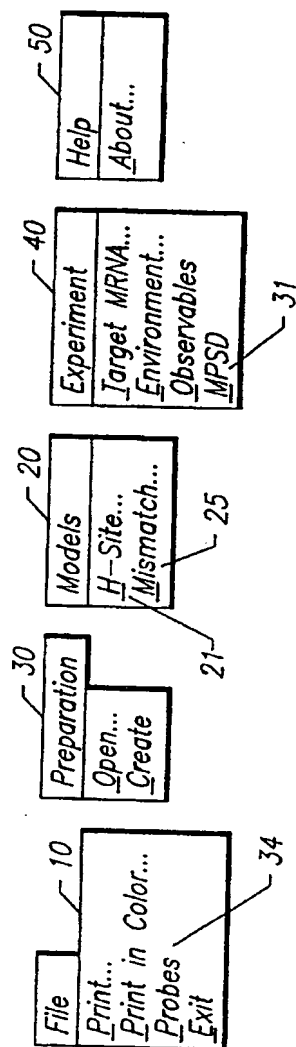
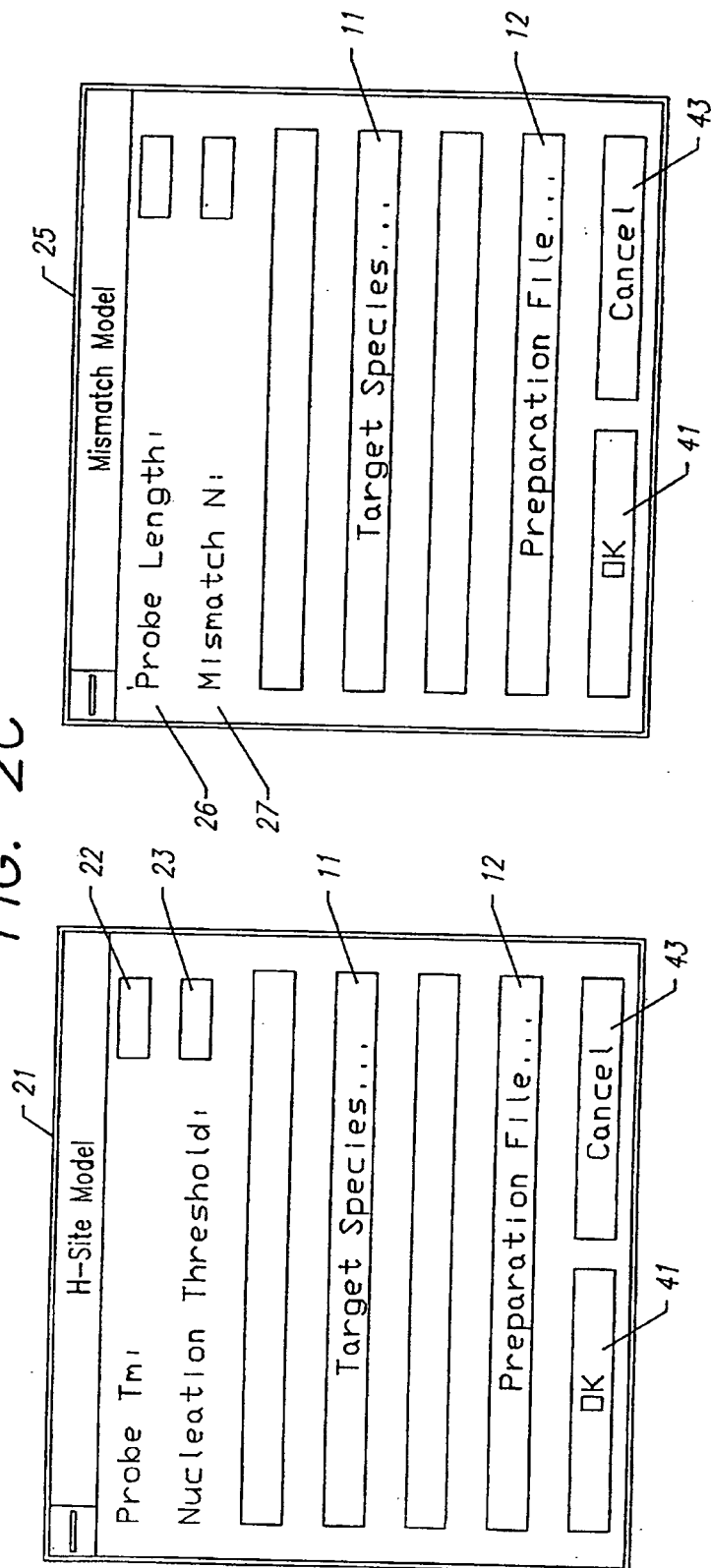


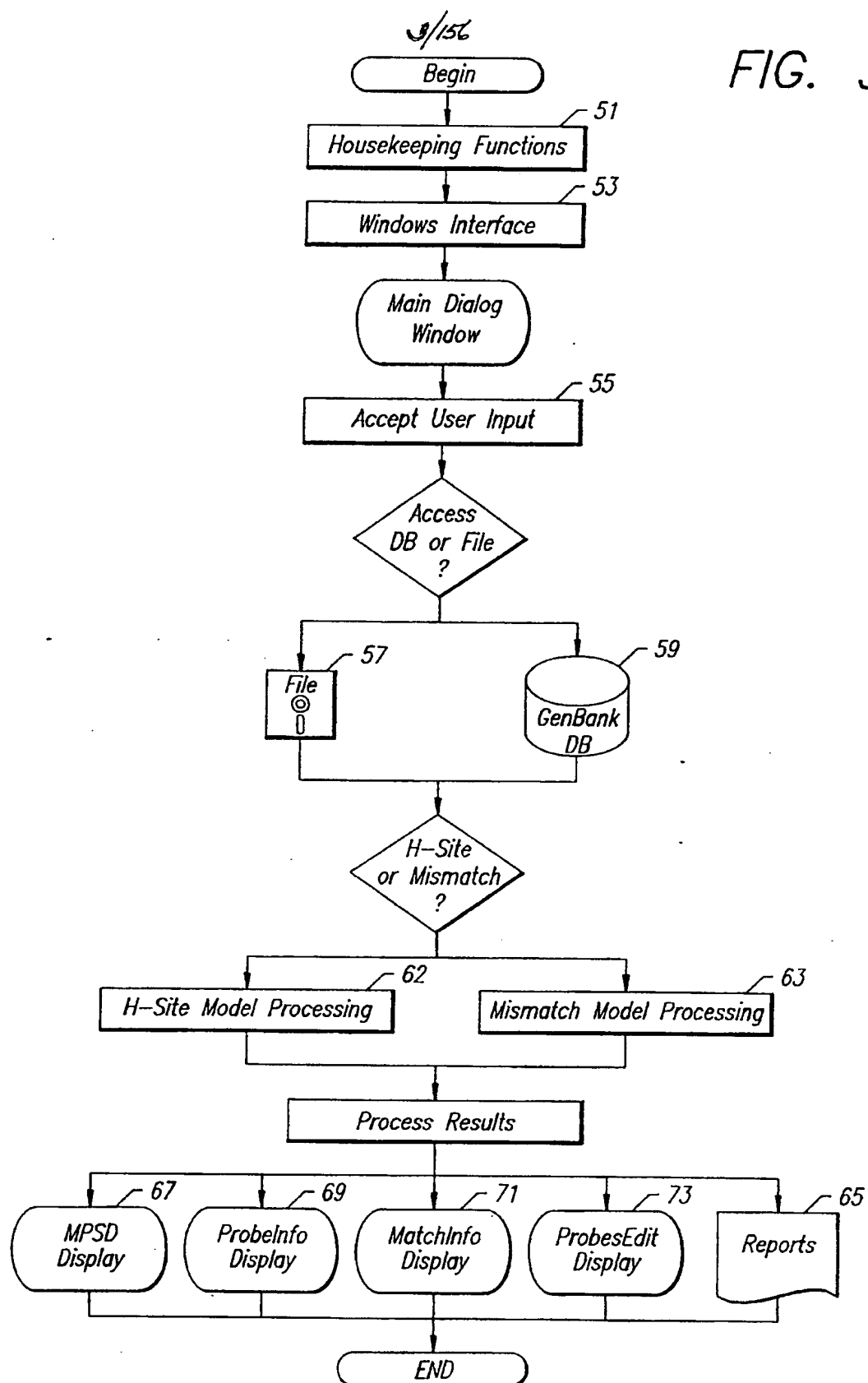
FIG. 2B

FIG. 2C



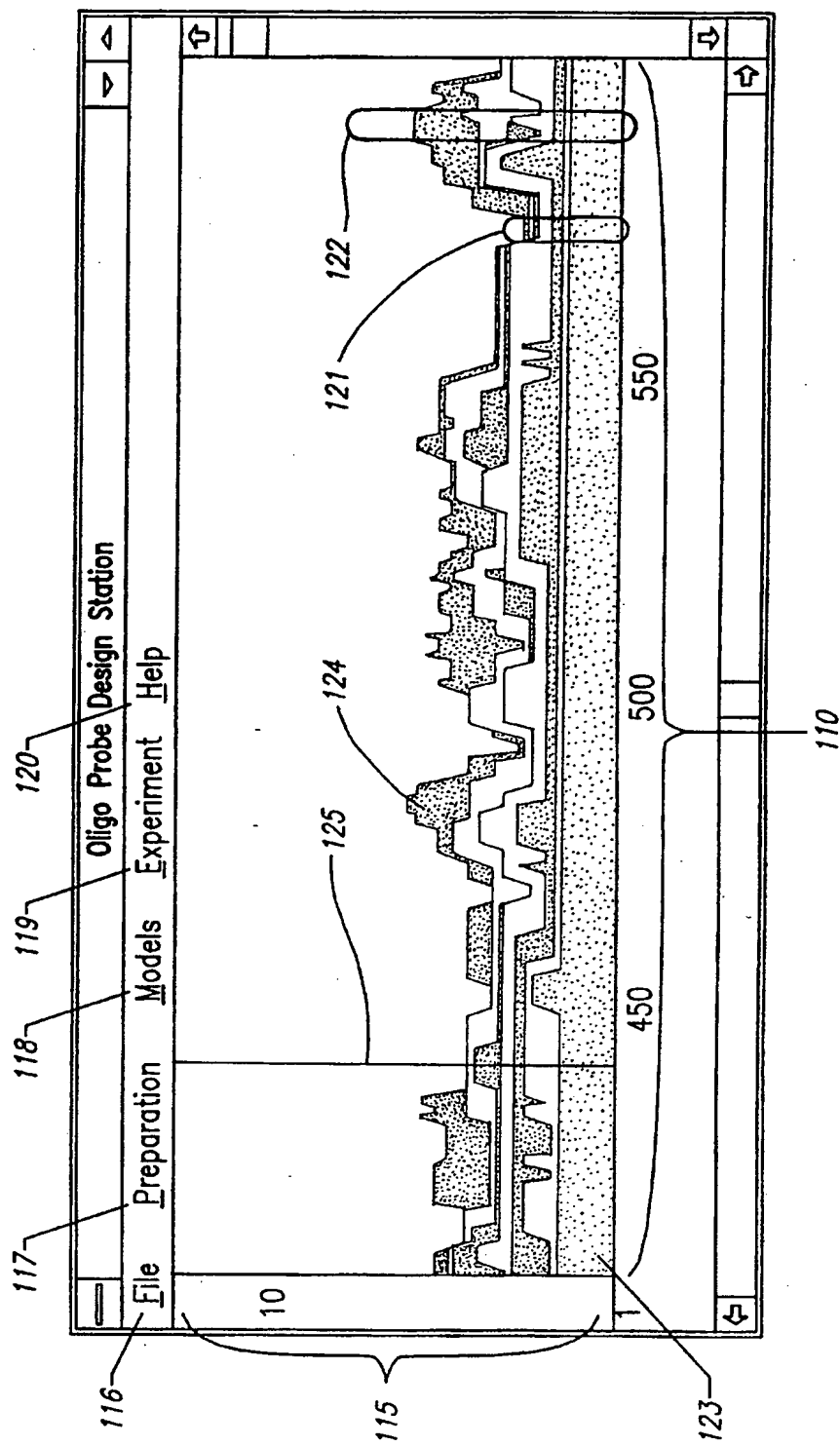
2/156

FIG. 3



4/156

FIG. 4



5/156

FIG. 5

Probe Info				
398	PROBE:	F: \MILAN\HUMBUNX.CDS		
	HYBRIDIZATIONS:	F: \MILAN\JUNMIX.PRP		
Locus	Pos	Tm	Length	HaIrpln
			21	4 1
humbjunx	398	0.0	agggcttcgcccga	cgctttg
muscjunx	398	61.7	-----	-----
humcjunx	323	50.0	-----	-----g-----c-
humdjunx	323	43.0	-----	-----g-----c-
humdjunx	215	36.2	ccctgc	-----gccc
humdjunx	401	36.0	---ag---	-----g-----c-
musdjunx	401	36.0	---ag---	-----a-----c-
humdjunx	100	35.7	g---gc---	-----cg-
musdjunx	262	34.3	ct-----	---gatct-gg-
humbjunx	659	30.5	c---cc---	-----gt---cacc
humdjunx	242	29.5	c---cacc---	-----c-gc
humdjunx	343	29.5	-cca-ca-	-----agc-ca
musbjunx	607	29.5	-ct-a-ac-	-----cacc
musdjunx	230	29.5	c-tgcg-	-----gccc
humcjunx	335	29.0	-----tgcg-	-----c-c-g-



6/156

FIG. 6

Probes selected - JUNMIX.prb	
File	
PROBE: C:\HITACHI\JUNMIX.PRP HYBRIDIZATION: C:\HITACHI\HUMBJU Length = 374 Hairpin = 3 5 Locus Pos Tm humbjunx 374 61.47 musbjunx 365 61.47 humdjunx 41 34.82 t-----9-9- humbjunx 182 31.12 a-----gt99 humdjunx 602 31.12 c-----c-99	155
PROBE: C:\HITACHI\JUNMIX.PRP HYBRIDIZATION: C:\HITACHI\HUMBJU Length = 467 Hairpin = 2 13 Locus Pos Tm humbjunx 467 61.7 musbjunx 458 51.6 humdjunx 32 29.35 tgagc99 humdjunx 32 29.35 tgagc99	156

7/156

## FIG. 6A (1)

PROBE: C:\HITACHI\JUNMIX.PRP  
 HYBRIDIZATION: C:\HITACHI\HUMBJUNX.CDS  
 Length = 374 Hairpin = 3 5  
 Locus      Pos   Tm  
 humbjunx 374 61.47 -----  
 musbjunx 365 61.47 -----  
 humdjunx 41  34.82 t-----g-g--agt  
 humbjunx 182 31.12 a-----gtgg--gc  
 humdjunx 602 31.12 c-----c-ggg-gc  
 humdjunx 602 31.12 c-----c-ggg-gc

PROBE: C:\HITACHI\JUNMIX.PRP  
 HYBRIDIZATION: C:\HITACHI\HUMBJUNX.CDS  
 Length = 377 Hairpin = 2 14  
 Locus      Pos   Tm  
 humbjunx 377 61.55 -----  
 musbjunx 368 61.55 -----  
 humdjunx 383 28.12 tg-cg-c--g-----  
 musdjunx 383 28.12 tg-ca-c--g-----  
 musdjunx 383 28.12 tg-ca-c--g-----

PROBE: C:\HITACHI\JUNMIX.PRP  
 HYBRIDIZATION: C:\HITACHI\HUMBJUNX.CDS  
 Length = 389 Hairpin = 3 3  
 Locus      Pos   Tm  
 humbjunx 389 61.7    -----  
 muscjunx 314 56.65 -c-----  
 musbjunx 380 50.85 -----t--g  
 humcjunx 314 49.35 -t-----g-----  
 humdjunx 395 33.85 -----tt-gc--ag  
 musdjunx 395 33.85 -----tt-gc--aa  
 humcjunx 326 32.35 g-ttcgcc-----tg  
 humdjunx 404 32.35 --ttcgcc-----t-  
 muscjunx 326 32.35 gcttcgcc-----tg  
 musdjunx 253 30.85 gacg-gct-ct-----  
 humbjunx 953 30.65 g-----t--c-cagct-  
 musdjunx 83  27.3    cc-gcggg-gt-----g

8/156

## FIG. 6A (2)

PROBE: C:\HITACHI\JUNMIX.PRP

HYBRIDIZATION: C:\HITACHI\HUMBJUNX.CDS

Length = 397 Hairpin = 4 1

Locus	Pos	Tm	
humbjunx	397	61.55	-----
muscjunx	322	53.44	-----g---
humcjunx	322	45.33	-----g-----g---
musbjunx	388	41.38	-----t--g-----t
humdjunx	214	36.83	ccccctgc-----
humdjunx	99	36.16	cg---gc-c-----
musdjunx	261	34.55	-ct-----gatct
humdjunx	400	33.27	c---ag-----g---
musdjunx	400	33.27	c---ag-----a---
humcjunx	334	32.28	-----tgcg--c-
humdjunx	412	32.28	-----t-a-g-c-
muscjunx	334	32.28	-----tgcg--c-
humbjunx	658	30.17	cc-cc-----gt---
humdjunx	241	28.95	-c--cacc-c-----
humdjunx	342	28.95	c-cca-ca-----ag
musbjunx	606	28.95	---ct-a-ac-----
musdjunx	229	28.95	-c-ctgcg-c-----
musdjunx	91	26.67	-gt-----gcc-ccg

PROBE: C:\HITACHI\JUNMIX.PRP

HYBRIDIZATION: C:\HITACHI\HUMBJUNX.CDS

Length = 417 Hairpin = 2 15

Locus	Pos	Tm	
humbjunx	417	60.08	-----
musbjunx	408	55.52	-----c-----
humdjunx	420	37.3	c-----g-----g---t-a-
musbjunx	61	29.0	g---gg-----ca-cctgt-
muscjunx	672	26.27	gc-gc-----a-g--aga--

9/156

## FIG. 6A (3)

PROBE: C:\HITACHI\JUNMIX.PRP  
 HYBRIDIZATION: C:\HITACHI\HUMBJUNX.CDS  
 Length = 461 Hairpin = 4 9  

Locus	Pos	Tm	
humbjunx	461	61.63	-----
musbjunx	452	61.63	-----
musbjunx	452	61.63	-----

PROBE: C:\HITACHI\JUNMIX.PRP  
 HYBRIDIZATION: C:\HITACHI\HUMBJUNX.CDS  
 Length = 467 Hairpin = 2 13  

Locus	Pos	Tm	
humbjunx	467	61.7	-----
musbjunx	458	51.6	-----c-g-
humdjunx	32	29.35	tgagcgg-----gcgg-
humdjunx	32	29.35	tgagcgg-----gcgg-

PROBE: C:\HITACHI\JUNMIX.PRP  
 HYBRIDIZATION: C:\HITACHI\HUMBJUNX.CDS  
 Length = 477 Hairpin = 2 4  

Locus	Pos	Tm	
humbjunx	477	61.37	-----
humdjunx	489	34.93	c-c---cg-----
humdjunx	489	34.93	c-c---cg-----

PROBE: C:\HITACHI\JUNMIX.PRP  
 HYBRIDIZATION: C:\HITACHI\HUMBJUNX.CDS  
 Length = 487 Hairpin = 3 3  

Locus	Pos	Tm	
humbjunx	487	61.14	-----
musdjunx	74	51.0	ct-----
humdjunx	499	45.64	-----t---g
humdjunx	527	30.72	cc-c-c-----
musdjunx	97	30.72	ttc-c-----g
musdjunx	580	30.72	-cc-----t-g
musdjunx	637	30.72	cc-cc-----g
musdjunx	637	30.72	cc-cc-----g

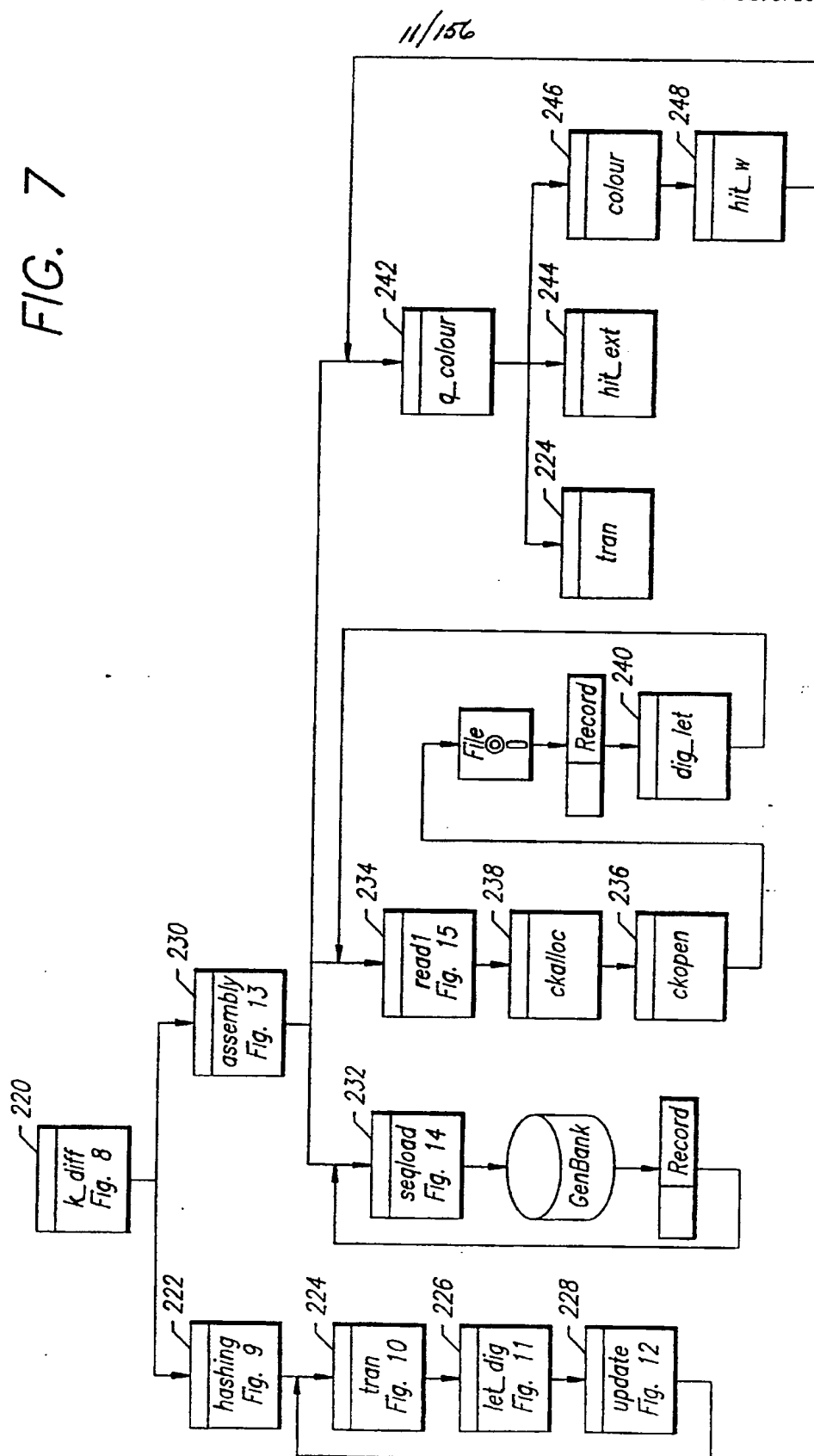
10/156

## FIG. 6A (4)

PROBE: C:\HITACHI\JUNMIX.PRP  
HYBRIDIZATION: C:\HITACHI\HUMBJUNX.CDS  
Length = 498 Hairpin = 3 2  
Locus Pos Tm  
humbjunx 498 61.26 -----  
humbjunx 498 61.26 -----

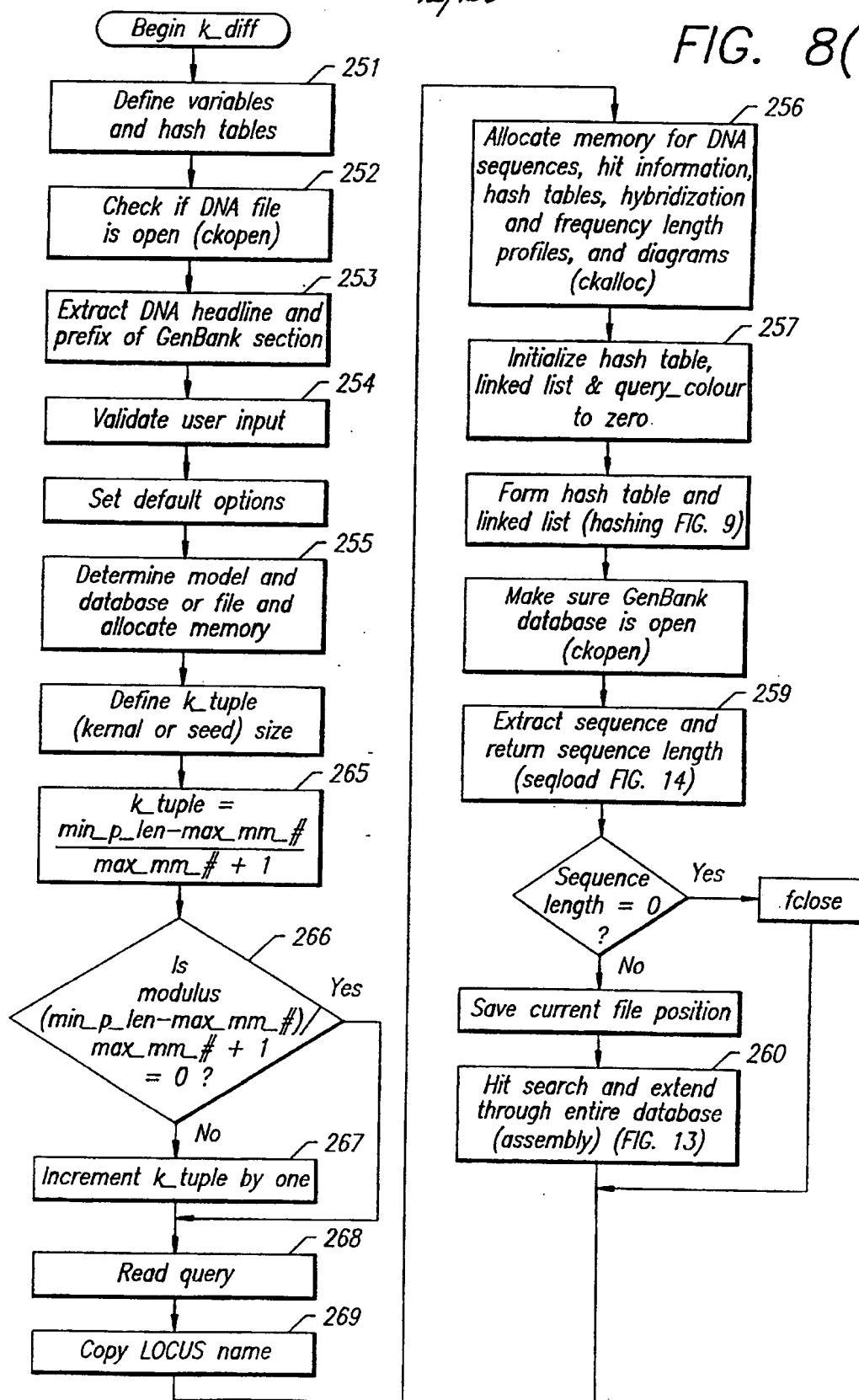
PROBE: C:\HITACHI\JUNMIX.PRP  
HYBRIDIZATION: C:\HITACHI\HUMBJUNX.CDS  
Length = 504 Hairpin = 3 2  
Locus Pos Tm  
humbjunx 504 61.47 -----  
musbjunx 495 40.35 c--a-----t-  
humdjunx 609 35.29 cg-----cgggg-  
humdjunx 609 35.29 cg-----cgggg-

FIG. 7



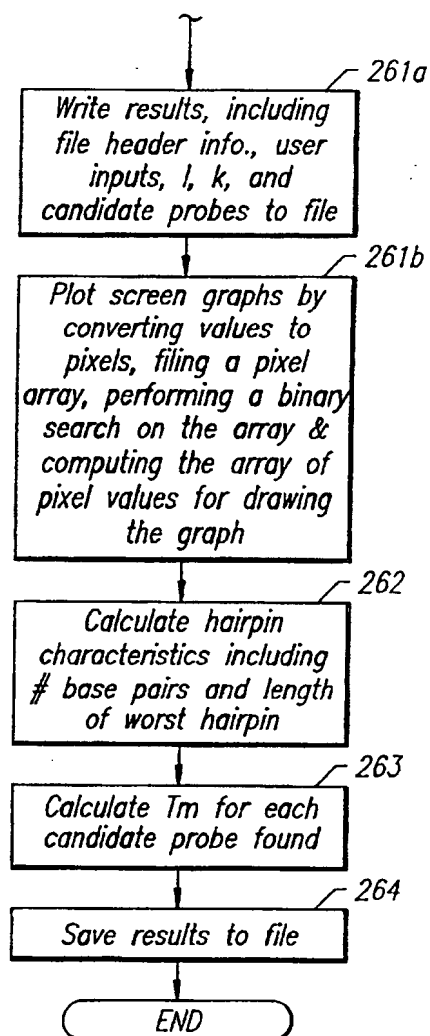
12/156

FIG. 8(1)



13/156

FIG. 8(2)





14/156

FIG. 9

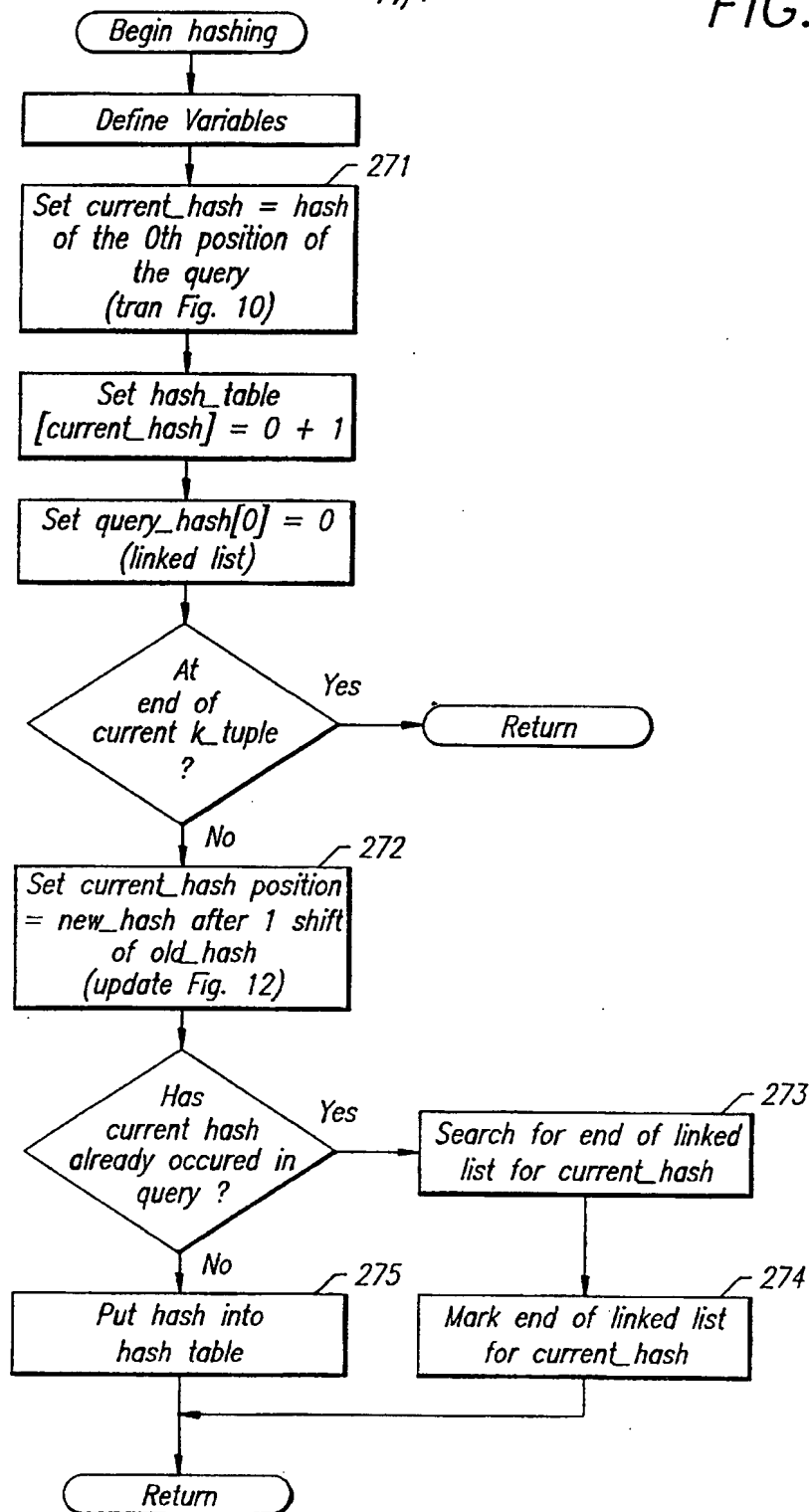


FIG. 10

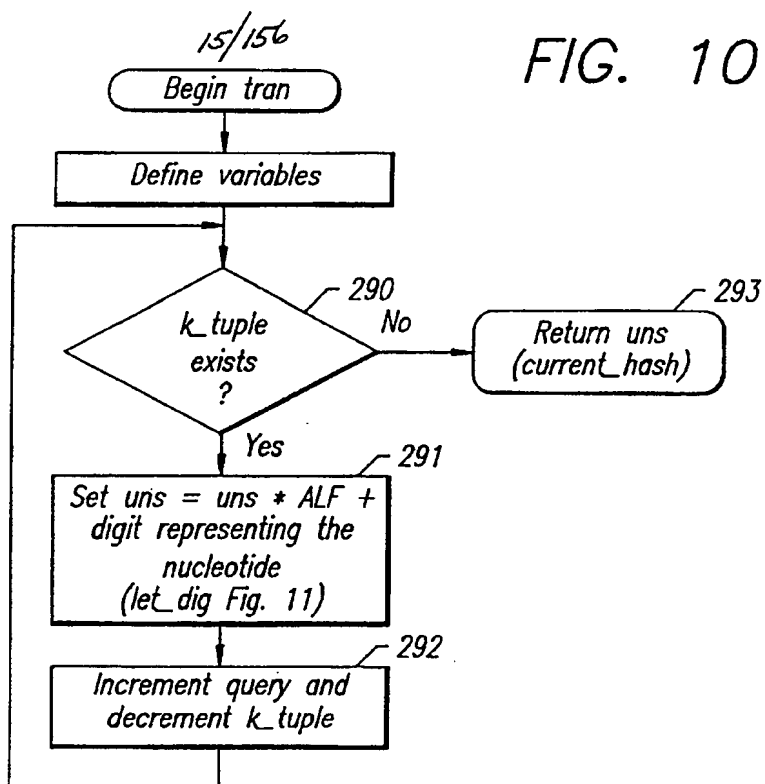
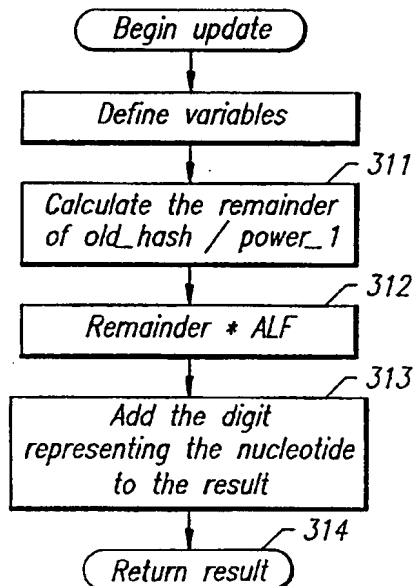


FIG. 12



16/156

FIG. 11

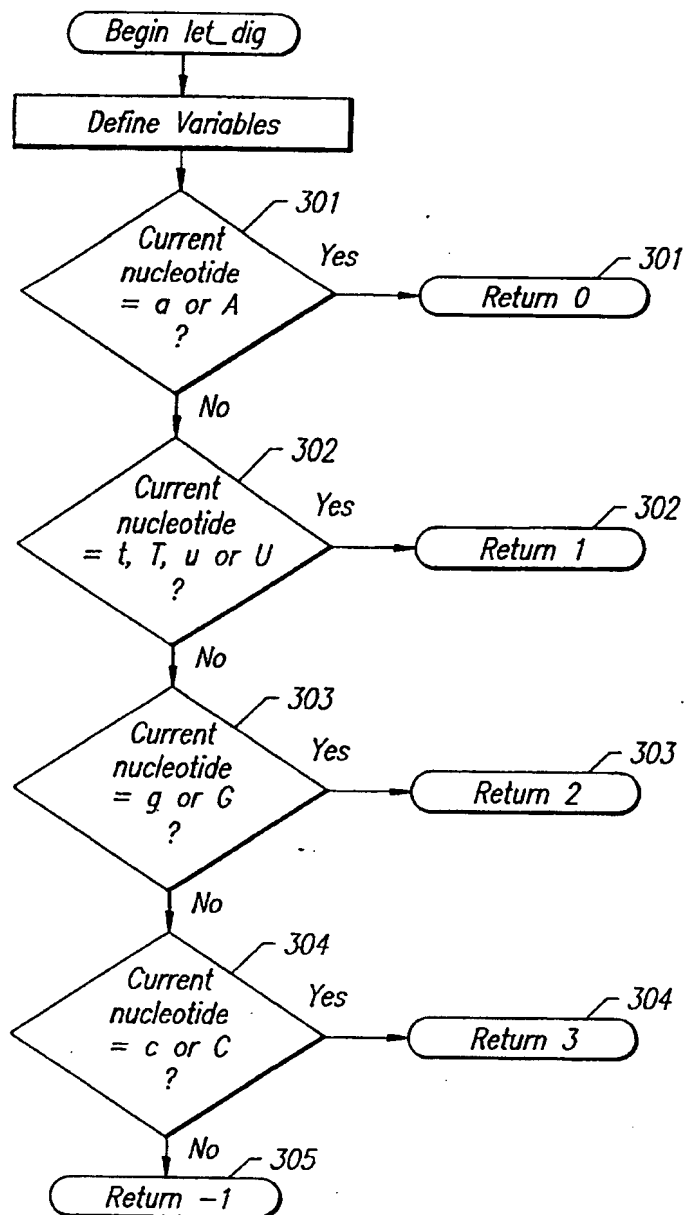
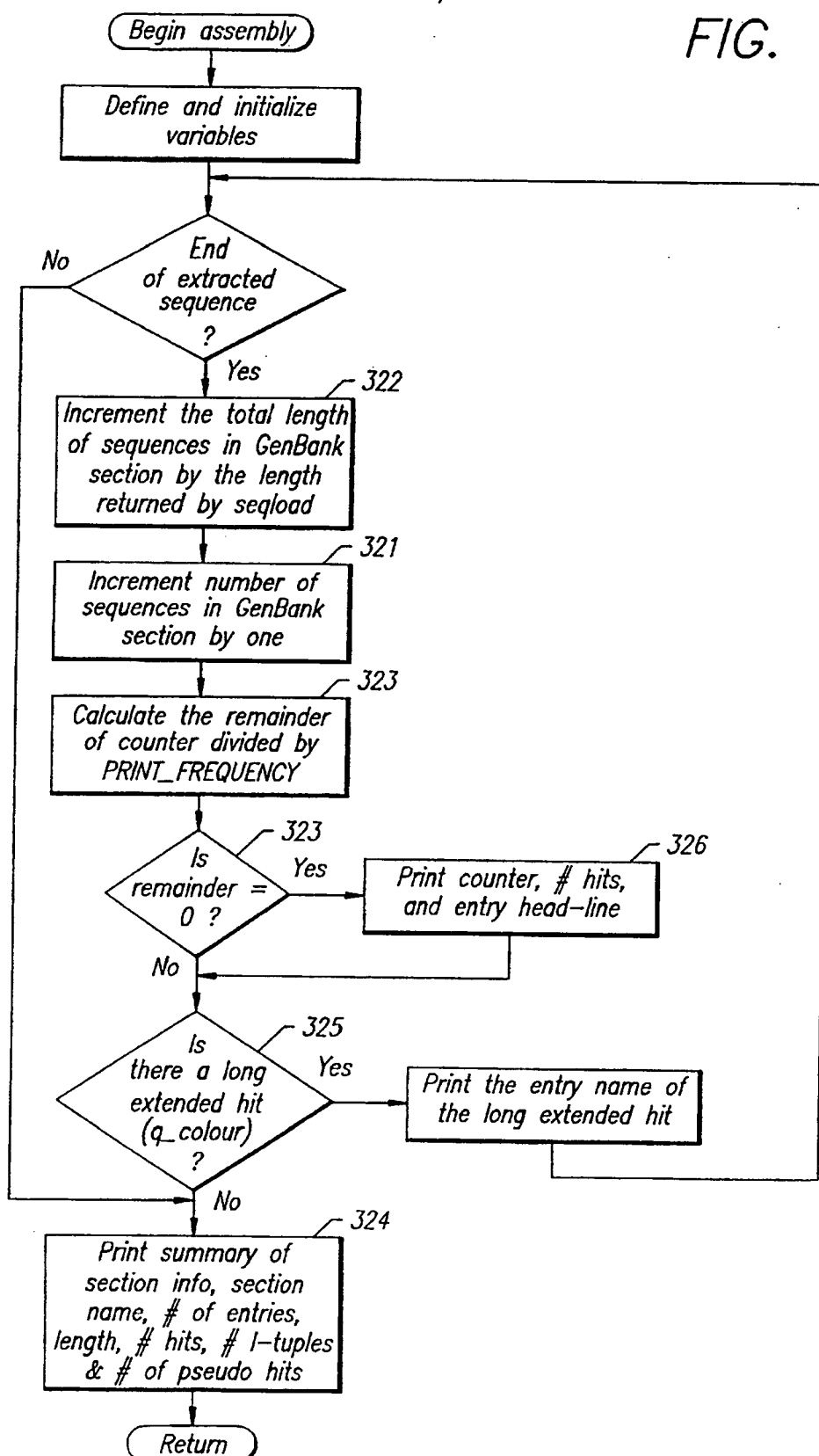
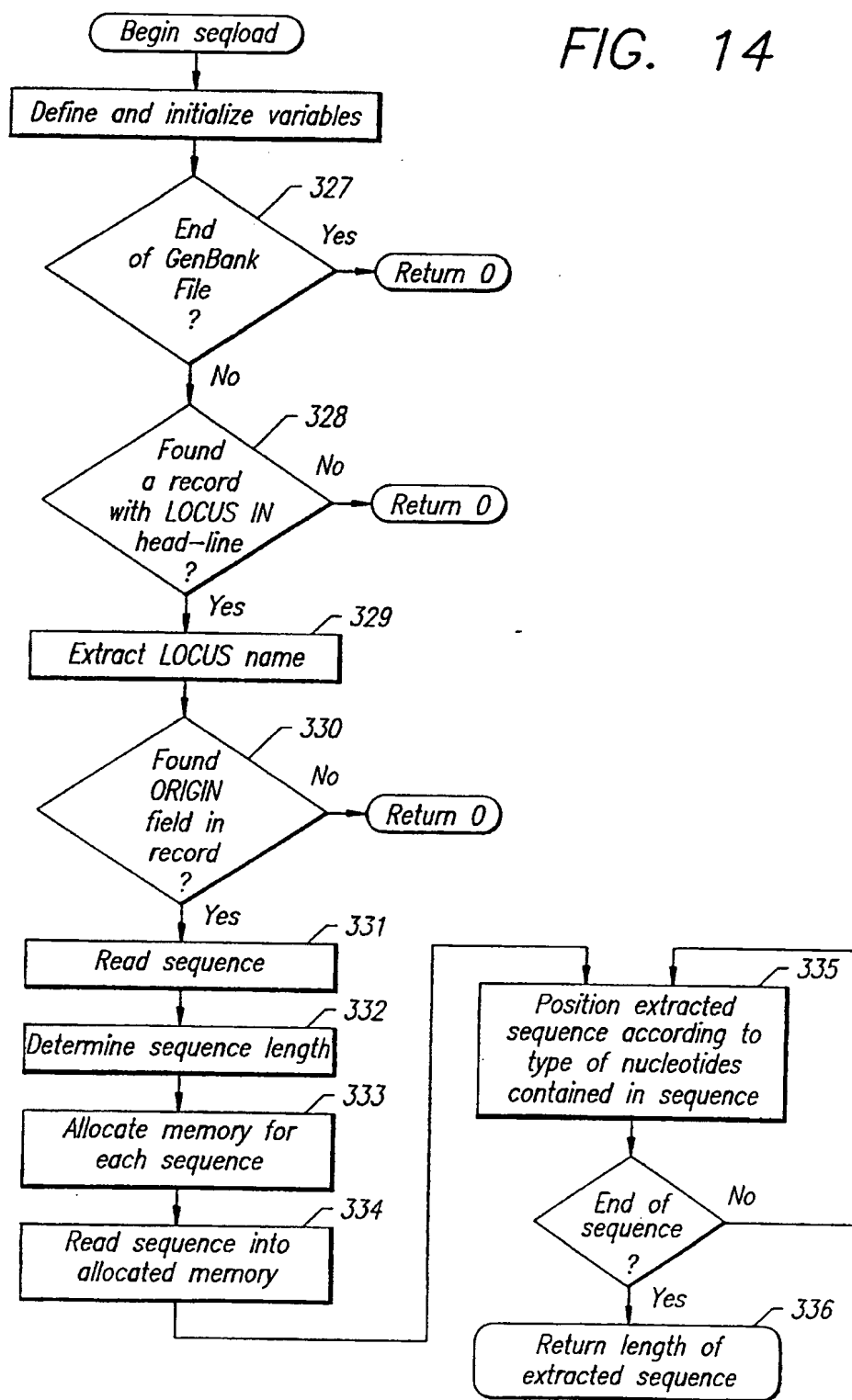


FIG. 13



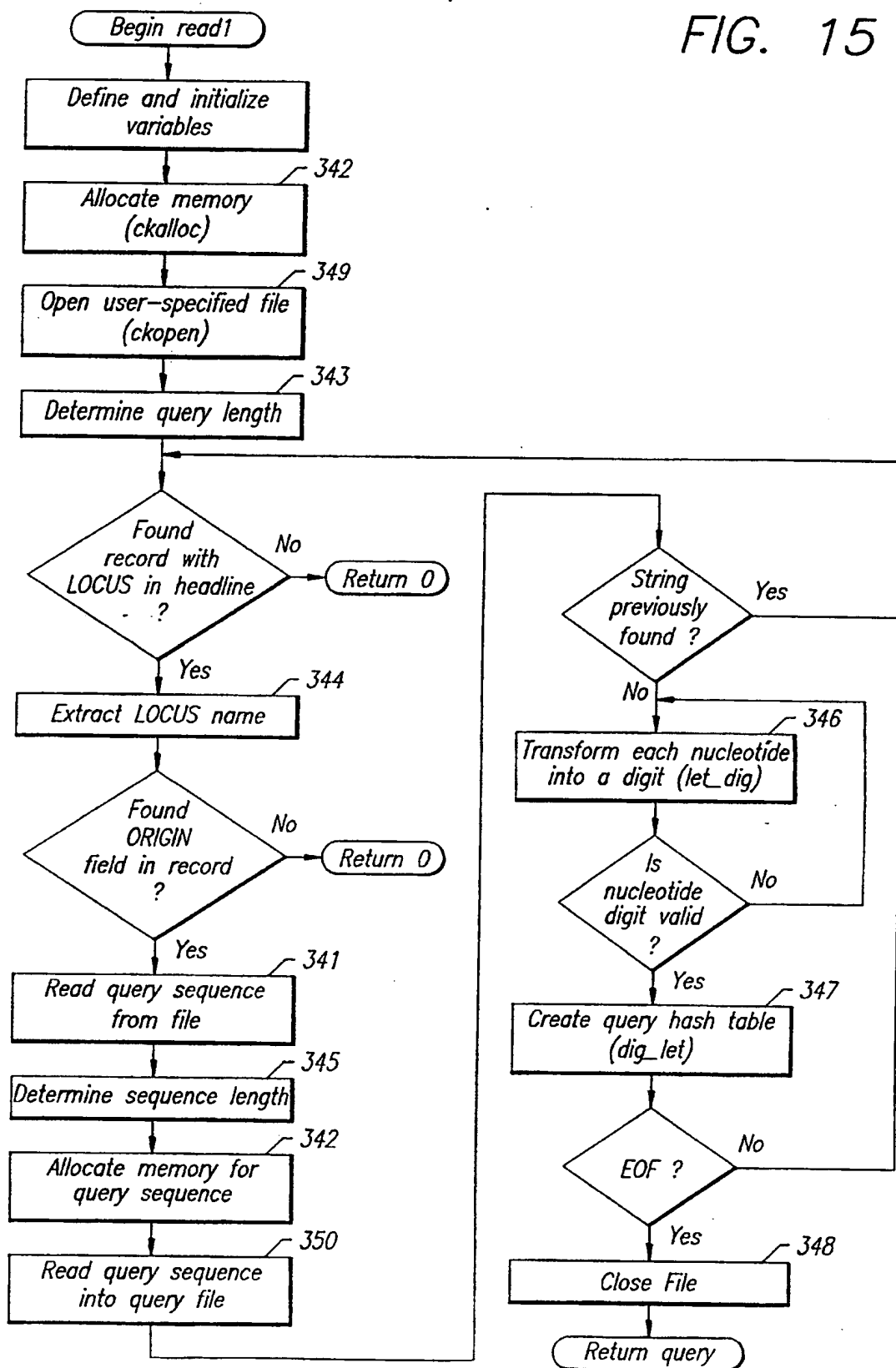
18/156

FIG. 14



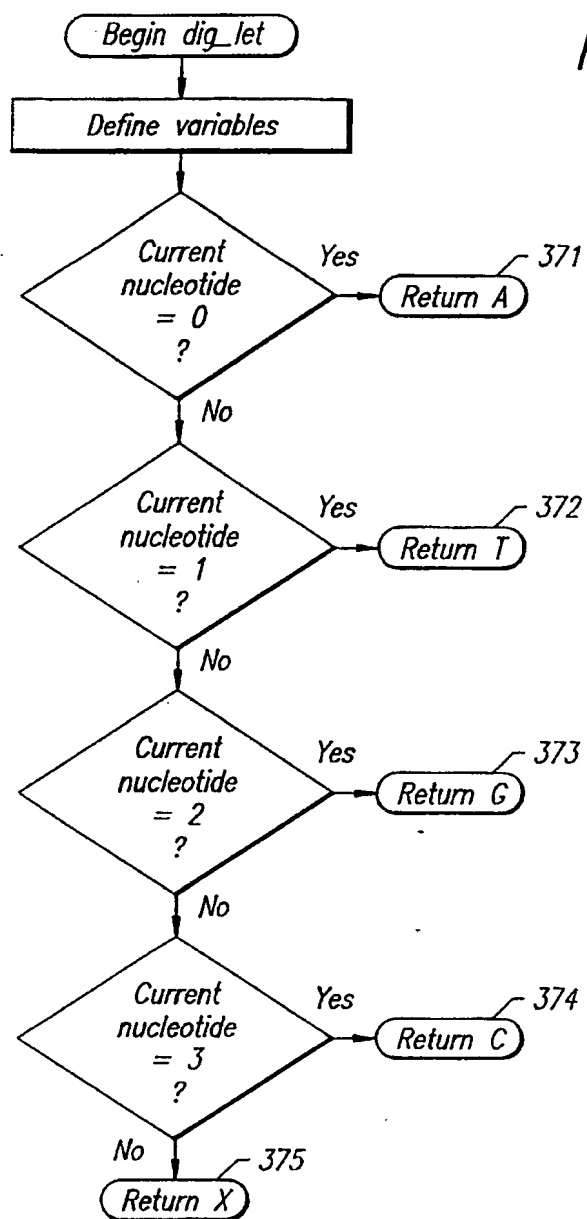
19/156

FIG. 15



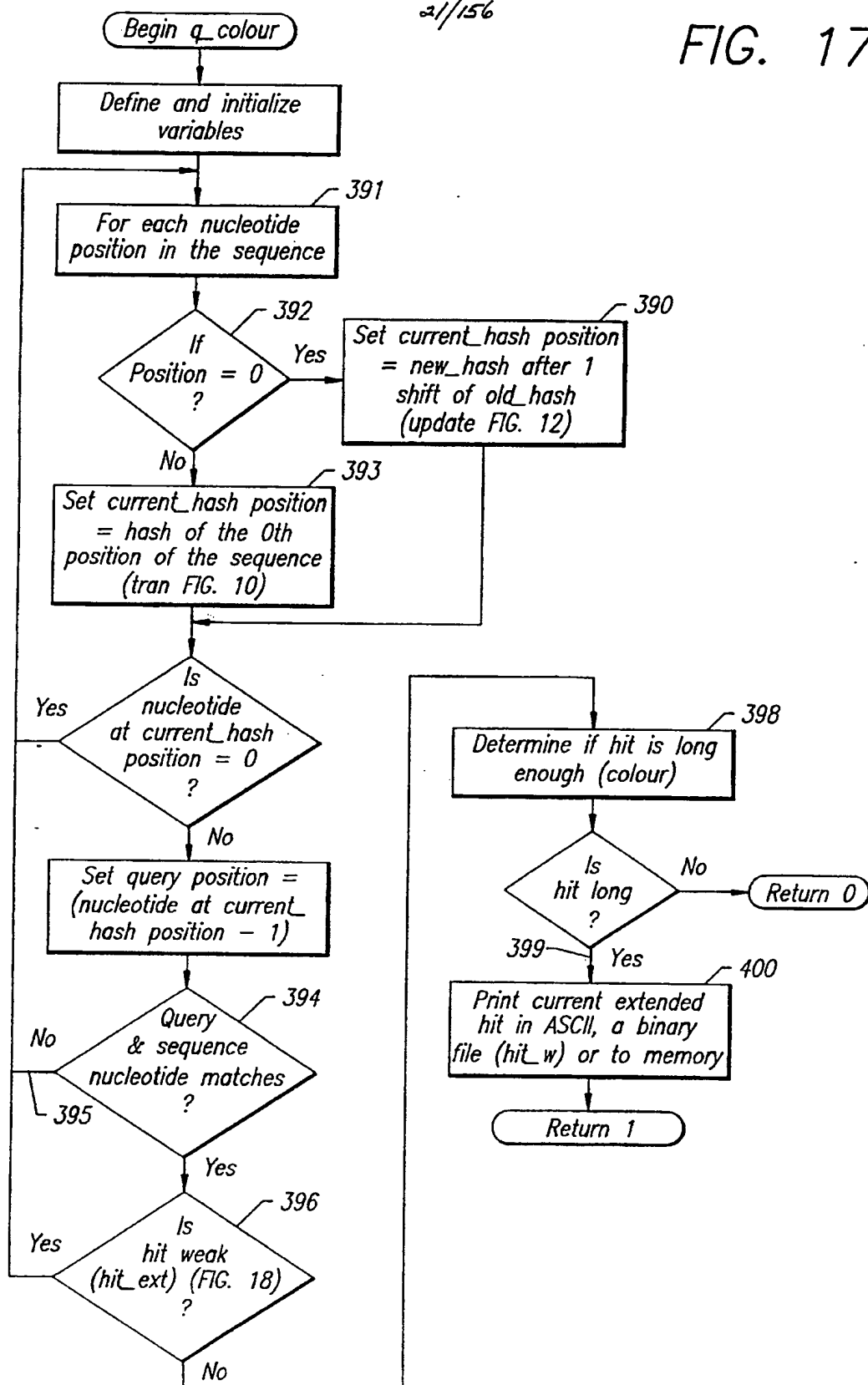
20/156

FIG. 16



21/156

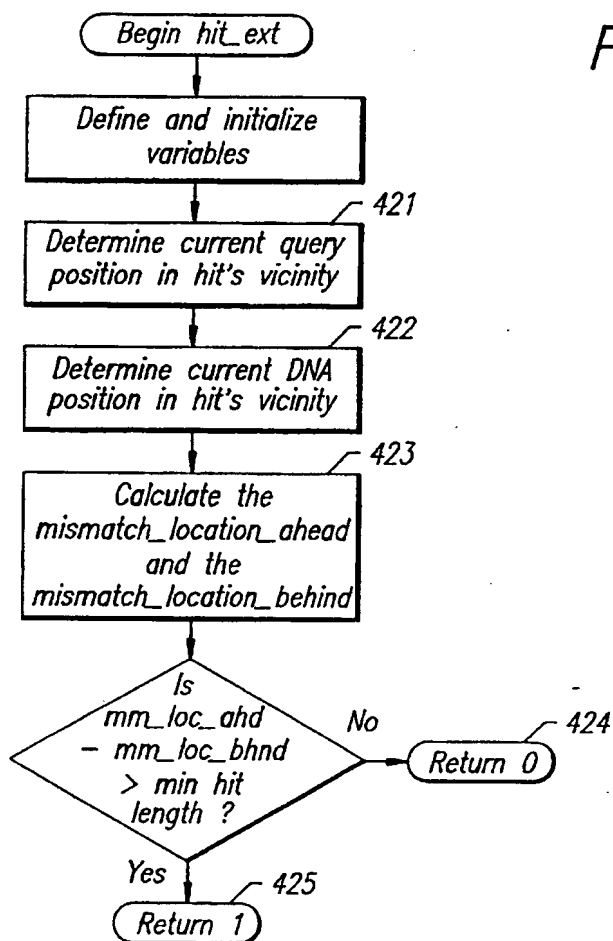
FIG. 17





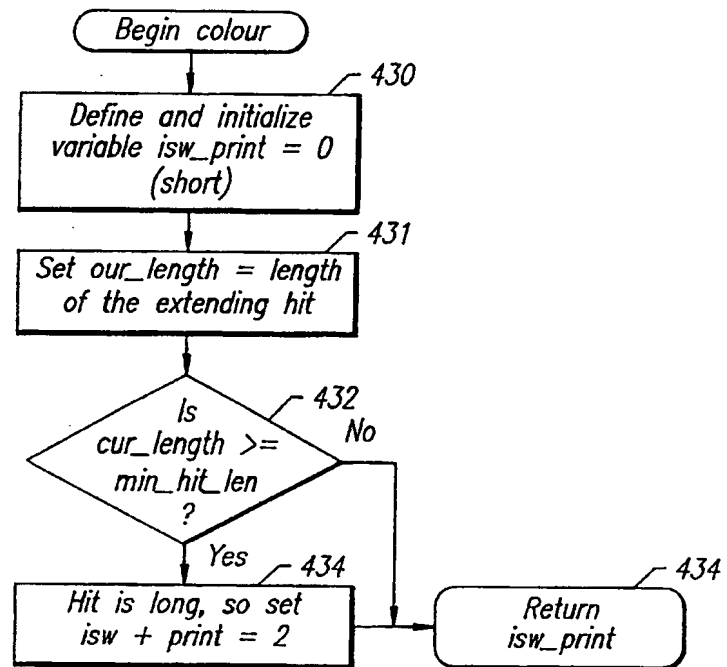
22/156

FIG. 18



23/156

FIG. 19



# FIG. 20 (1)

OligoProbe DesignStation

Probes: C:\HITACHI\HUMBJUNX.CDS  
 Database: C:\HITACHI\JUNMIX.SEQ

Mismatch Model, 1 = 21, k = 4

24/156

Position length	Mismatches				screensN				Probe
	0	1	2	3	4	5	6	7	8
1 21	1	1	1	1	1	ATGTGCACTAAAAATGGAACAG			
2 21	1	1	1	1	1	TGTGCACTAAAAATGGAACAGC			
3 21	1	1	1	1	1	GTGCACTAAAAATGGAACAGCC			
4 21	1	1	1	1	1	TGCACTAAAAATGGAACAGCCC			
5 21	1	1	1	1	1	GCACTAAAAATGGAACAGCCCT			
6 21	1	1	1	1	1	CACTAAAAATGGAACAGCCCTT			
7 21	1	1	1	1	1	ACTAAAAATGGAACAGCCCTTC			
8 21	1	1	1	1	1	CTAAAAATGGAACAGCCCTTCT			
9 21	1	1	1	1	1	TAAAAATGGAACAGCCCTTCTA			
10 21	1	1	1	1	1	AAAAATGGAACAGCCCTTCTAC			
11 21	1	1	1	1	1	AAATGGAACAGCCCTTCTACC			
12 21	1	1	1	1	1	AATGGAACAGCCCTTCTACCA			
13 21	1	1	1	1	1	ATGGAACAGCCCTTCTACCAC			
14 21	1	1	1	1	1	TGGAACAGCCCTTCTACCACG			



26/156

FIG. 20 (3)

37	21	1	1	1	1	1	GACTCATACACAGCTACGGGA
38	21	1	1	1	1	1	ACTCATACACAGCTACGGGAT
39	21	1	1	1	1	1	CTCATACACAGCTACGGGATA
40	21	1	1	1	1	1	TCATACACAGCTACGGGATAC
41	21	1	1	1	1	1	CATACACAGCTACGGGATACG
42	21	1	1	1	1	1	ATACACAGCTACGGGATACGG
43	21	1	1	1	1	1	TACACAGCTACGGGATACGGC
44	21	1	1	1	1	1	ACACAGCTACGGGATACGGCC
45	21	1	1	1	1	1	CACAGCTACGGGATACGGCCG
46	21	1	1	1	1	1	ACAGCTACGGGATACGGCCGG
47	21	1	1	1	1	1	CAGCTACGGGATACGGCCGGG
48	21	1	1	1	1	1	AGCTACGGGATACGGCCGGGC
49	21	1	1	1	1	1	GCTACGGGATACGGCCGGGCC
50	21	1	1	1	1	1	CTACGGGATACGGCCGGGCC
51	21	1	1	1	1	1	TACGGGATACGGCCGGGCC
52	21	1	1	1	1	1	ACGGGATACGGCCGGGCCCT
53	21	1	1	1	1	1	CGGGATACGGCCGGGCCCTG

27/156

FIG. 20 (4)

54	21	1	1	1	1	1	GGGATACGGCCGGGCCCCCTGG
55	21	1	1	1	1	1	GGATACGGCCGGGCCCCCTGGT
56	21	1	1	1	1	1	GATACGGCCGGGCCCCCTGGTG
57	21	1	1	1	1	1	ATACGGCCGGGCCCCCTGGTGG
58	21	1	1	1	1	1	TACGGCCGGGCCCCCTGGTGGC
59	21	1	1	1	1	1	ACGGCCGGGCCCCCTGGTGGCC
60	21	1	1	1	1	1	CGGCCGGGCCCCCTGGTGGCCT
61	21	1	1	1	1	1	GGCCGGGCCCCCTGGTGGCCTC
62	21	1	1	1	1	1	GCCGGGCCCCCTGGTGGCCTCT
63	21	1	1	1	1	1	CCGGCCCTGGTGGCCTCTC
64	21	1	1	1	1	1	CGGCCCTGGTGGCCTCTCT
65	21	1	1	1	1	1	GGCCCCCTGGTGGCCTCTCTC
66	21	1	1	1	1	1	GGCCCTGGTGGCCTCTCTCT
67	21	1	1	1	1	1	GCCCCCTGGTGGCCTCTCTCTA
68	21	1	1	1	1	1	CCCTGGTGGCCTCTCTCTAC
69	21	1	1	1	1	1	CCCTGGTGGCCTCTCTCTACA
70	21	1	1	1	1	1	CCTGGTGGCCTCTCTCTACAC
71	21	1	1	1	1	1	CTGGTGGCCTCTCTCTACACG
72	21	1	1	1	1	1	TGGTGGCCTCTCTCTACACGA
73	21	1	1	1	1	1	GGTGGCCTCTCTCTACACGAC
74	21	1	1	1	1	1	GTGGCCTCTCTCTACACGACT

28/156

FIG. 20 (5)

75	1	1	1	1	1	TGGCCTCTCTCTACACGACTA
76	1	1	1	1	1	GGCCTCTCTCTACACGACTAC
77	1	1	1	1	1	GCCTCTCTCTACACGACTACA
78	1	1	1	1	1	CCTCTCTCTACACGACTACAA
79	1	1	1	1	1	CTCTCTCTACACGACTACAAA
80	1	1	1	1	1	TCTCTCTACACGACTACAAAC
81	1	1	1	1	1	CTCTCTACACGACTACAAACT
82	1	1	1	1	1	TCTCTACACGACTACAAACTC
83	1	1	1	1	1	CTCTACACGACTACAAACTCC
84	1	1	1	1	1	TCTACACGACTACAAACTCCT
85	1	1	1	1	1	CTACACGACTACAAACTCCTG
86	1	1	1	1	1	TACACGACTACAAACTCCTGA
87	1	1	1	1	1	ACACGACTACAAACTCCTGAA
88	1	1	1	1	1	CACGACTACAAACTCCTGAAA
89	1	1	1	1	1	ACGACTACAAACTCCTGAAAC
90	1	1	1	1	1	CGACTACAAACTCCTGAAACC
91	1	1	1	1	1	GACTACAAACTCCTGAAACCG
92	1	1	1	1	1	ACTACAAACTCCTGAAACCGA
93	1	1	1	1	1	CTACAAACTCCTGAAACCGAG
94	1	1	1	1	1	TACAAACTCCTGAAACCGAGC
95	1	1	1	1	1	ACAAACTCCTGAAACCGAGCC









32/156

FIG. 20 (9)

[illegible]

33/156

FIG. 20 (10)

181	21	1	1	1	1	1	GAGGGCGGCGGTGGCGGCAGC
182	21	1	1	1	1	1	AGGGCGGCGGTGGCGGCAGCT
183	21	1	1	1	1	1	GGCGGCGGTTGGCGGCAGCTA
184	21	1	1	1	1	1	GGCGGCGGTGGCGGCAGCTAC
185	21	1	1	1	1	1	GCGGCGGTGGCGGCAGCTACT
186	21	1	1	1	1	1	CGGCGGTGGCGGCAGCTACTT
187	21	1	1	1	1	1	GGCGGTGGCGGCAGCTACTTT
188	21	1	1	1	1	1	GCGTGCGGCGCAGCTACTTTT
189	21	1	1	1	1	1	CGGTGGCGGCAGCTACTTTTC
190	21	1	1	1	1	1	GGTGGCGGCAGCTACTTTTCT
191	21	1	1	1	1	1	GTGGCGGCAGCTACTTTTCTG
192	21	1	1	1	1	1	TGGCGGCAGCTACTTTTCTGG
193	21	1	1	1	1	1	GGCGGCAGCTACTTTTCTGGT
194	21	1	1	1	1	1	GCGGCAGCTACTTTTCTGGTC
195	21	1	1	1	1	1	CGGCAGCTACTTTTCTGGTCA
196	21	1	1	1	1	1	GGCAGCTACTTTTCTGGTCAG
197	21	1	1	1	1	1	GCAGCTACTTTTCTGGTCAGG
198	21	1	1	1	1	1	CAGCTACTTTTCTGGTCAGGG
199	21	1	1	1	1	1	AGCTACTTTTCTGGTCAGGGC
200	21	1	1	1	1	1	GCTACTTTTCTGGTCAGGGCT
201	21	1	1	1	1	1	CTACTTTTCTGGTCAGGGCTC

34/156

FIG. 20 (11)

202	21	1	1	1	1	1	TACTTTTCTGGTCAGGGCTCG
203	21	1	1	1	1	1	ACTTTTCTGGTCAGGGCTCGG
204	21	1	1	1	1	1	CTTTTCTGGTCAGGGCTCGGA
205	21	1	1	1	1	1	TTTTCCTGGTCAGGGCTCGGAC
206	21	1	1	1	1	1	TTTCTGGTCAGGGCTCGGACA
207	21	1	1	1	1	1	TTCTGGTCAGGGCTCGGACAC
208	21	1	1	1	1	1	TCTGGTCAGGGCTCGGACACC
209	21	1	1	1	1	1	CTGGTCAGGGCTCGGACACCG
210	21	1	1	1	1	1	TGGTCAGGGCTCGGACACCGG
211	21	1	1	1	1	1	GGTCAGGGCTCGGACACCGGC
212	21	1	1	1	1	1	GTCAGGGCTCGGACACCGGCG
213	21	1	1	1	1	1	TCAGGGCTCGGACACCGGCGC
214	21	1	1	1	1	1	CAGGGCTCGGACACCGGCGCG
215	21	1	1	1	1	1	AGGGCTCGGACACCGGCGCGT
216	21	1	1	1	1	1	GGGCTCGGACACCGGCGCGTC
217	21	1	1	1	1	1	GGCTCGGACACCGGCGCGTCT
218	21	1	1	1	1	1	GCTCGGACACCGGCGCGTCTC
219	21	1	1	1	1	1	CTCGGACACCGGCGCGTCTCT
220	21	1	1	1	1	1	TCGGACACCGGCGCGTCTCTC
221	21	1	1	1	1	1	CGGACACCGGCGCGTCTCTCA
222	21	1	1	1	1	1	GGACACCGGCGCGTCTCTCAA



36/56

FIG. 20 (13)

245	1	1	1	1	1	TCGCCTCTTTCGGAGCTGGAAC
246	1	1	1	1	1	CGCCTCTTTCGGAGCTGGAACG
247	1	1	1	1	1	GCCTCTTTCGGAGCTGGAACGC
248	1	1	1	1	1	CCTCTTTCGGAGCTGGAACGCC
249	1	1	1	1	1	CTCTTTCGGAGCTGGAACGCCT
250	1	1	1	1	1	TCTTCGGAGCTGGAACGCCCTG
251	1	1	1	1	1	CTTCGGAGCTGGAACGCCCTGA
252	1	1	1	1	1	TTCGGAGCTGGAACGCCCTGAT
253	1	1	1	1	1	TCGGAGCTGGAACGCCCTGATT
254	1	1	1	1	1	CGGAGCTGGAACGCCCTGATTG
255	1	1	1	1	1	GGAGCTGGAACGCCCTGATTGT
256	1	1	1	1	1	GAGCTGGAACGCCCTGATTGTC
257	1	1	1	1	1	AGCTGGAACGCCCTGATTGTCC
258	1	1	1	1	1	GCTGGAACGCCCTGATTGTCCC
259	1	1	1	1	1	CTGGAACGCCCTGATTGTCCCC
260	1	1	1	1	1	TGGAACGCCCTGATTGTCCCCA
261	1	1	1	1	1	GGAACGCCCTGATTGTCCCCAA
262	1	1	1	1	1	GAACGCCCTGATTGTCCCCAAC
263	1	1	1	1	1	AACGCCCTGATTGTCCCCAACA
264	1	1	1	1	1	ACGCCCTGATTGTCCCCAACAG
265	1	1	1	1	1	CGCCTGATTGTCCCCAACAGC

37/156

FIG. 20 (14)

266	21	1	1	1	1	1	GCCTGATTGTCCCCAACAGCA	1
267	21	1	1	1	1	1	CCTGATTGTCCCCAACAGCAA	1
268	21	1	1	1	1	1	CTGATTGTCCCCAACAGCAAC	1
269	21	1	1	1	1	1	TGATTGTCCCCAACAGCAACG	1
270	21	1	1	1	1	1	GATTGTCCCCAACAGCAACGG	1
271	21	1	1	1	1	1	ATTGTCCCCAACAGCAACGGC	1
272	21	1	1	1	1	1	TTGTCCCCAACAGCAACGGCG	1
273	21	1	1	1	1	1	TGTCCCCAACAGCAACGGCGT	1
274	21	1	1	1	1	1	GTCCCCAACAGCAACGGCGTG	1
275	21	1	1	1	1	1	TCCCCAACAGCAACGGCGTGA	1
276	21	1	1	1	1	1	CCCCAACAGCAACGGCGTGAT	1
277	21	1	1	1	1	1	CCCAACAGCAACGGCGTGATC	1
278	21	1	1	1	1	1	CCAACAGCAACGGCGTGATCA	1
279	21	1	1	1	1	1	CAACAGCAACGGCGTGATCAC	1
280	21	1	1	1	1	1	AACAGCAACGGCGTGATCAG	1
281	21	1	1	1	1	1	ACAGCAACGGCGTGATCACGA	1
282	21	1	1	1	1	1	CAGCAACGGCGTGATCACCAG	1
283	21	1	1	1	1	1	AGCAACGGCGTGATCACCAGG	1
284	21	1	1	1	1	1	GCAACGGCGTGATCACCAGCA	1
285	21	1	1	1	1	1	CAACGGCGTGATCACCAGCAC	1
286	21	1	1	1	1	1	AACGGCGTGATCACCAGCAG	1



38/156

[illegible]

39/156

FIG. 20 (16)

[illegible]

40/152

FIG. 20 (17)

330	21	1	1	1	1	1	TTACCCCGCGGGGTGGCAG	1
331	21	1	1	1	1	1	TACCCCGCGGGGTGGCAGC	1
332	21	1	1	1	1	1	ACCCCGCGGGGTGGCAGCG	1
333	21	1	1	1	1	1	CCCCCGCGGGGTGGCAGCGG	1
334	21	1	1	1	1	1	CCCCCGCGGGGTGGCAGCGGT	1
335	21	1	1	1	1	1	CCCGCGGGGTGGCAGCGGTG	1
336	21	1	1	1	1	1	CCCGCGGGGTGGCAGCGGTGG	1
337	21	1	1	1	1	1	CGCGGGGTGGCAGCGGTGGA	1
338	21	1	1	1	1	1	GCGGGGTGGCAGCGGTGGAG	1
339	21	1	1	1	1	1	CGGGGTGGCAGCGGTGGAGG	1
340	21	1	1	1	1	1	GGGGTGGCAGCGGTGGAGGT	1
341	21	1	1	1	1	1	GGGTGGCAGCGGTGGAGGTG	1
342	21	1	1	1	1	1	GGTGGCAGCGGTGGAGGTGC	1
343	21	1	1	1	1	1	GTGGCAGCGGTGGAGGTGCA	1
344	21	1	1	1	1	1	GTGGCAGCGGTGGAGGTGCAG	1
345	21	1	1	1	1	1	TGGCAGCGGTGGAGGTGCAGG	1
346	21	1	1	1	1	1	GGCAGCGGTGGAGGTGCAGGG	1
347	21	1	1	1	1	1	GCAGCGGTGGAGGTGCAGGGG	1
348	21	1	1	1	1	1	CAGCGGTGGAGGTGCAGGGGG	1
349	21	1	1	1	1	1	AGCGGTGGAGGTGCAGGGGGC	1
350	21	1	1	1	1	1	GCGGTGGAGGTGCAGGGGGCG	1

41/156

FIG. 20 (18)

351	21	1	1	1	1	1	CGGTGGAGGTGCAGGGGGCGC
352	21	1	1	1	1	1	GGTGGAGGTGCAGGGGGCGCA
353	21	1	1	1	1	1	GTGAGGTGCAGGGGGCGCAG
354	21	1	1	1	1	1	TGGAGGTGCAGGGGGCGCAGG
355	21	1	1	1	1	1	GGAGGTGCAGGGGGCGCAGGG
356	21	1	1	1	1	1	GAGTGCAGGGGGCGCAGGGG
357	21	1	1	1	1	1	AGGTGCAGGGGGCGCAGGGG
358	21	1	1	1	1	1	GGTGCAGGGGGCGCAGGGGG
359	21	1	1	1	1	1	GTGAGGGGGCGCAGGGGGCG
360	21	1	1	1	1	1	TGCAGGGGGCGCAGGGGGCGG
361	21	1	1	1	1	1	GCAGGGGGCGCAGGGGGCGGC
362	21	1	1	1	1	1	CAGGGGGCGCAGGGGGCGGCG
363	21	1	1	1	1	1	AGGGGGCGCAGGGGGCGGCGT
364	21	1	1	1	1	1	GGGGGGCGCAGGGGGCGGCGTC
365	21	1	1	1	1	1	GGGGGGCGCAGGGGGCGGCGTCA
366	21	1	1	1	1	1	GGGGGGCGCAGGGGGCGGCGTCAC
367	21	1	1	1	1	1	GGGGGGCGCAGGGGGCGGCGTCACC
368	21	1	1	1	1	1	GGGGGGCGCAGGGGGCGGCGTCACCG
369	21	1	1	1	1	1	GGGGGGCGCAGGGGGCGGCGTCACCGA
370	21	1	1	2	2	2	GGGGGGCGCAGGGGGCGGCGTCACCGAG
371	21	1	2	2	2	2	GGGGGGCGCAGGGGGCGGCGTCACCGAGG
372	21	1	2	2	2	2	GGGGGGCGCAGGGGGCGGCGTCACCGAGGA

42/156

FIG. 20 (19)

[illegible]

43/156

FIG. 20 (20)

394	21	1	2	2	2	2	2	2	CAGGAGGGCTTCGCCGACGGC
395	21	1	2	2	2	2	2	2	AGGAGGGCTTCGCCGACGGCT
396	21	1	2	2	2	2	2	2	GGAGGGCTTCGCCGACGGCTT
397	21	1	1	2	1	1	1	1	GAGGGCTTCGCCGACGGCTTT
398	21	1	1	1	1	1	1	1	AGGGCTTCGCCGACGGCTTTG
399	21	1	1	1	1	1	1	1	GGCTTCGCCGACGGCTTTGT
400	21	1	1	1	1	1	1	1	GGCTTCGCCGACGGCTTTGTC
401	21	1	1	1	1	1	1	1	GCTTCGCCGACGGCTTTGTCA
402	21	1	1	1	1	1	1	1	CTTCGCCGACGGCTTTGTCAA
403	21	1	1	1	1	1	1	1	TTCGCCGACGGCTTTGTCAAA
404	21	1	1	1	1	1	1	1	TCGCCGACGGCTTTGTCAAAG
405	21	1	1	1	1	1	1	1	CGCCGACGGCTTTGTCAAAGC
406	21	1	1	1	1	1	1	1	GCCGACGGCTTTGTCAAAGCC
407	21	1	1	1	1	1	1	1	CCGACGGCTTTGTCAAAGCCC
408	21	1	1	1	1	1	1	1	CGACGGCTTTGTCAAAGCCCT
409	21	1	2	2	2	2	2	2	GACGGCTTTGTCAAAGCCCCTG
410	21	1	2	2	2	2	2	2	ACGGCTTTGTCAAAGCCCCTGG
411	21	1	2	2	2	2	2	2	CGGCTTTGTCAAAGCCCCTGGA
412	21	1	2	2	2	2	2	2	GGCTTTGTCAAAGCCCCTGGAC
413	21	1	2	2	2	2	2	2	GCTTTGTCAAAGCCCCTGGACG
414	21	1	2	2	2	2	2	2	CTTTGTCAAAGCCCCTGGACGA

44/152

FIG. 20 (21)

415	21	1	2	2	2	2	TTTGTCAAAGCCCTGGACGAT
416	21	1	2	2	2	2	TTGTCAAAGCCCTGGACGATC
417	21	1	2	2	2	2	TGTCAAAGCCCTGGACGATCT
418	21	1	2	2	2	2	GTCAAAGCCCTGGACGATCTG
419	21	1	2	2	2	2	TCAAAGCCCTGGACGATCTGC
420	21	1	2	2	2	2	CAAAGCCCTGGACGATCTGCA
421	21	1	2	2	2	2	AAAGCCCTGGACGATCTGCAC
422	21	1	2	2	2	2	AAGCCCTGGACGATCTGCACA
423	21	1	2	2	2	2	AGCCCTGGACGATCTGCACAA
424	21	1	2	2	2	2	GCCCTGGACGATCTGCACAAG
425	21	1	2	2	2	2	CCCTGGACGATCTGCACAAGA
426	21	1	2	2	2	2	CCTGGACGATCTGCACAAGAT
427	21	1	2	2	2	2	CTGGACGATCTGCACAAGATG
428	21	1	2	2	2	2	TGGACGATCTGCACAAGATGA
429	21	1	2	2	2	2	GGACGATCTGCACAAGATGAA
430	21	1	2	2	2	2	GACGATCTGCACAAGATGAAC
431	21	1	2	2	2	2	ACGATCTGCACAAGATGAACC
432	21	1	2	2	2	2	CGATCTGCACAAGATGAACCA
433	21	1	2	2	2	2	GATCTGCACAAGATGAACCCAC
434	21	1	2	2	2	2	ATCTGCACAAGATGAACCCACG
435	21	1	2	2	2	2	TCTGCACAAGATGAACCCACGT
436	21	2	2	2	2	2	CTGCACAAGATGAACCCACGTG







47/156

FIG. 20 (24)

[illegible]

48/156

FIG. 20 (25)

501	21	1	1	1	1	1	TGGGCCCGGGGGCGTCTACGC	1
502	21	1	1	1	1	1	GGCCCCGGGGCGTCTACGCC	1
503	21	1	1	1	1	1	GGCCCGGGGGCGTCTACGCCG	1
504	21	1	1	1	1	1	GCCCCGGGGCGTCTACGCCGG	1
505	21	1	1	1	1	1	CCCGGGGGCGTCTACGCCGGC	1
506	21	1	1	1	1	1	CCGGGGCGTCTACGCCGGCC	1
507	21	1	1	1	1	1	CGGGGGCGTCTACGCCGGCCC	1
508	21	1	1	1	1	1	GGGGCGTCTACGCCGGCCCCG	1
509	21	1	1	1	1	1	GGGGCGTCTACGCCGGCCCCG	1
510	21	1	1	1	1	1	GGCGTCTACGCCGGCCCCGA	1
511	21	1	1	1	1	1	GGCGTCTACGCCGGCCCCGAG	1
512	21	1	1	1	1	1	GCGTCTACGCCGGCCCCGAGC	1
513	21	1	1	1	1	1	CGTCTACGCCGGCCCCGAGCC	1
514	21	1	1	1	1	1	GTCTACGCCGGCCCCGAGCCA	1
515	21	1	1	1	1	1	TCTACGCCGGCCCCGAGCCAC	1
516	21	1	1	1	1	1	CTACGCCGGCCCCGAGCCACC	1
517	21	1	1	1	1	1	TACGCCGGCCCCGAGCCACCT	1
518	21	1	1	1	1	1	ACGCCGGCCCCGAGCCACCTC	1
519	21	1	1	1	1	1	CGCCGGCCCCGAGCCACCTCC	1
520	21	1	1	1	1	1	GCCGGCCCCGAGCCACCTCCC	1
521	21	1	1	1	1	1	CCGGCCCCGAGCCACCTCCCG	1

49/156

FIG. 20 (26)

[illegible]

50/156

FIG. 20 (27)

543	21	1	1	1	1	1	TACACCAACCTCAGCAGCTA
544	21	1	2	2	2	2	TACACCAACCTCAGCAGCTAC
545	21	1	2	2	2	2	ACACCAACCTCAGCAGCTACT
546	21	1	2	2	2	2	CACCAACCTCAGCAGCTACTC
547	21	1	1	1	1	1	ACCAACCTCAGCAGCTACTCC
548	21	1	1	1	1	1	CCAACCTCAGCAGCTACTCCC
549	21	1	1	1	1	1	CAACCTCAGCAGCTACTCCCC
550	21	1	1	1	1	1	AACCTCAGCAGCTACTCCCCA
551	21	1	1	1	1	1	ACCTCAGCAGCTACTCCCCAG
552	21	1	1	1	1	1	CCTCAGCAGCTACTCCCCCAGC
553	21	1	1	1	1	1	CTCAGCAGCTACTCCCCCAGCC
554	21	1	1	1	1	1	TCAGCAGCTACTCCCCCAGCCT
555	21	1	1	1	1	1	CAGCAGCTACTCCCCCAGCCTC
556	21	1	1	1	1	1	AGCAGCTACTCCCCCAGCCTCT
557	21	1	1	1	1	1	GCAGTACTCCCCCAGCCTCTG
558	21	1	1	1	1	1	CAGTACTCCCCCAGCCTCTGCG
559	21	1	1	1	1	1	AGTACTCCCCCAGCCTCTGCG
560	21	1	1	1	1	1	GCTACTCCCCCAGCCTCTGCCGT
561	21	1	1	1	1	1	CTACTCCCCCAGCCTCTGCCGTC
562	21	1	1	1	1	1	TACTCCCCCAGCCTCTGCCGTCC
563	21	1	1	1	1	1	ACTCCCCCAGCCTCTGCCGT CCT
564	21	1	1	1	1	1	CTCCCCCAGCCTCTGCCGT CCTC

51/156

FIG. 20 (28)

565	1	1	1	1	1	TCCCCAGCCTCTGCGTCCTCG
566	1	1	1	1	1	CCCCAGCCTCTGCGTCCTCGG
567	1	1	1	1	1	CCCAGCCTCTGCGTCCTCGGG
568	1	1	1	1	1	CCAGCCTCTGCGTCCTCGGGA
569	1	1	1	1	1	CAGCCTCTGCGTCCTCGGGAG
570	1	1	1	1	1	AGCCTCTGCGTCCTCGGGAGG
571	1	1	1	1	1	GCCTCTGCGTCCTCGGGAGGC
572	1	1	1	1	1	CCTCTGCGTCCTCGGGAGGCG
573	1	1	1	1	1	CTCTGCGTCCTCGGGAGGCGC
574	1	1	1	1	1	TCTGCGTCCTCGGGAGGCGCC
575	1	1	1	1	1	CTGCGTCCTCGGGAGGCGCCG
576	1	1	1	1	1	TGCGTCCTCGGGAGGCGCCGG
577	1	1	1	1	1	GCGTCCTCGGGAGGCGCCGGG
578	1	1	1	1	1	CGTCCTCGGGAGGCGCCGGGG
579	1	1	1	1	1	GTCTCGGGAGGCGCCGGGGC
580	1	1	1	1	1	TCCTCGGGAGGCGCCGGGGCT
581	1	1	1	1	1	CCTCGGGAGGCGCCGGGGCTG
582	1	1	1	1	1	CTCGGGAGGCGCCGGGGCTGC
583	1	1	1	1	1	TCGGGAGGCGCCGGGGCTGCC
584	1	1	1	1	1	CGGAGGCGCCGGGGCTGCCG
585	1	1	1	1	1	GGGAGGCGCCGGGGCTGCCGT

52/156

FIG. 20 (29)

586	21	1	1	1	1	1	GGAGGCGCCGGGGCTGCCGTC
587	21	1	1	1	1	1	GAGGCGCCGGGGCTGCCGTCG
588	21	1	1	1	1	1	AGGCGCCGGGGCTGCCGTCGG
589	21	1	1	1	1	1	GGCGCCGGGGCTGCCGTCGGG
590	21	1	1	1	1	1	GCGCCGGGGCTGCCGTCGGGA
591	21	1	1	1	1	1	CGCCGGGGCTGCCGTCGGGAC
592	21	1	1	1	1	1	GCCGGGGCTGCCGTCGGGACC
593	21	1	1	1	1	1	CCGGGGCTGCCGTCGGGACCG
594	21	1	1	1	1	1	CGGGGCTGCCGTCGGGACCCG
595	21	1	1	1	1	1	GGGCTGCCGTCGGGACCCGGG
596	21	1	1	1	1	1	GGCTGCCGTCGGGACCCGGGA
597	21	1	1	1	1	1	GGCTGCCGTCGGGACCCGGGAG
598	21	1	1	1	1	1	GCTGCCGTCGGGACCCGGGAGC
599	21	1	1	1	1	1	CTGCCGTCGGGACCCGGGAGCT
600	21	1	1	1	1	1	TGCCGTCGGGACCCGGGAGCTC
601	21	1	1	1	1	1	GCCGTCGGGACCCGGGAGCTCG
602	21	1	1	1	1	1	CCGTCGGGACCCGGGAGCTCGT
603	21	1	1	1	1	1	CGTCGGGACCCGGGAGCTCGTA
604	21	1	1	1	1	1	GTCGGGACCCGGGAGCTCGTAC
605	21	1	1	1	1	1	TCGGGACCCGGGAGCTCGTACC
606	21	1	1	1	1	1	CGGGACCCGGGAGCTCGTACCC

53/156

607	21	1	1	1	1	1	1	1	GGACCGGAGCTCGTACCCG
608	21	1	1	1	1	1	1	1	GGACCGGAGCTCGTACCCGA
609	21	1	1	1	1	1	1	1	GACCGGAGCTCGTACCCGAC
610	21	1	1	1	1	1	1	1	ACCGGAGCTCGTACCCGACG
611	21	1	1	1	1	1	1	1	CCGGAGCTCGTACCCGACGA
612	21	1	1	1	1	1	1	1	C GGAGCTCGTACCCGACGAC
613	21	1	1	1	1	1	1	1	GGAGCTCGTACCCGACGACC
614	21	1	1	1	1	1	1	1	GGAGCTCGTACCCGACGACCA
615	21	1	1	1	1	1	1	1	GAGTCGTACCCGACGACCAC
616	21	1	1	1	1	1	1	1	AGTCGTACCCGACGACCACC
617	21	1	1	1	1	1	1	1	GCTCGTACCCGACGACCACCA
618	21	1	1	1	1	1	1	1	CTCGTACCCGACGACCACCAT
619	21	1	1	1	1	1	1	1	TCGTACCCGACGACCACCATC
620	21	1	1	1	1	1	1	1	CGTACCCGACGACCACCATCA
621	21	1	1	1	1	1	1	1	GTACCCGACGACCACCATCAG
622	21	1	1	2	2	2	2	2	TACCCGACGACCACCATCAGC
623	21	1	2	2	2	2	2	2	ACCCGACGACCACCATCAGCT
624	21	1	2	2	2	2	2	2	CCCGACGACCACCATCAGCTA
625	21	1	2	2	2	2	2	2	CCGACGACCACCATCAGCTAC
626	21	1	2	2	2	2	2	2	CGACGACCACCATCAGCTACC
627	21	1	2	2	2	2	2	2	GACGACCACCATCAGCTACCT
628	21	1	2	2	2	2	2	2	ACGACCACCATCAGCTACCTC





55/156

[illegible]

56/156

FIG. 20 (33)

[illegible]

57/156

FIG. 20 (34)

693	21	1	1	1	1	1	CTTGGGCGCGGCGCCTCCAC
694	21	1	1	1	1	1	TTGGGCGCGGCGCCTCCACC
695	21	1	1	1	1	1	TGGCGCGCGGCGCCTCCACCT
696	21	1	1	1	1	1	GGCGCGCGGCGCCTCCACCTT
697	21	1	1	1	1	1	GGCGCGCGGCGCCTCCACCTTC
698	21	1	1	1	1	1	GCGCGCGCGCCTCCACCTTCA
699	21	1	1	1	1	1	CCGCGCGCGCCTCCACCTTCAA
700	21	1	1	1	1	1	CGCGCGCGCCTCCACCTTCAAG
701	21	1	1	1	1	1	GCGCGCGCCTCCACCTTCAAGG
702	21	1	1	1	1	1	CGCGCGCCTCCACCTTCAAGGA
703	21	1	1	1	1	1	GCGCGCCTCCACCTTCAAGGAG
704	21	1	1	1	1	1	GCGCCTCCACCTTCAAGGAGG
705	21	1	1	1	1	1	CGCCTCCACCTTCAAGGAGGA
706	21	1	1	1	1	1	GCCTCCACCTTCAAGGAGGAA
707	21	1	1	1	1	1	CCTCCACCTTCAAGGAGGAAC
708	21	1	1	1	1	1	CTCCACCTTCAAGGAGGAACC
709	21	1	1	1	1	1	TCCACCTTCAAGGAGGAACCG
710	21	1	1	1	1	1	CCACCTTCAAGGAGGAACCGC
711	21	1	1	1	1	1	CACCTTCAAGGAGGAACCGCA
712	21	1	1	1	1	1	ACCTTCAAGGAGGAACCGCAG
713	21	1	1	1	1	1	CCTTCAAGGAGGAACCGCAGA

58/156

FIG. 20 (35)

714	21	1	1	1	1	1	CTTCAAGGAGGAACCGCAGAC
715	21	1	1	1	1	1	TTCAAGGAGGAACCGCAGACC
716	21	1	1	1	1	1	TCAAGGAGGAACCGCAGACCG
717	21	1	1	1	1	1	CAAGGAGGAACCGCAGACCGT
718	21	1	1	1	1	1	AAGGAGGAACCGCAGACCGTG
719	21	1	1	1	1	1	AGGAGGAACCGCAGACCGTGC
720	21	1	1	1	1	1	GGAGGAACCGCAGACCGTGCC
721	21	1	1	2	2	2	GAGGAACCGCAGACCGTGCCG
722	21	1	1	2	2	2	AGGAACCGCAGACCGTGCCGG
723	21	1	1	2	2	2	GGAACCGCAGACCGTGCCGGA
724	21	1	1	3	3	3	GAACCGCAGACCGTGCCGGAG
725	21	1	1	2	2	2	AACCGCAGACCGTGCCGGAGG
726	21	1	1	2	2	2	ACCGCAGACCGTGCCGGAGGC
727	21	1	1	1	1	1	CCGCAGACCGTGCCGGAGGCG
728	21	1	1	1	1	1	CGCAGACCGTGCCGGAGGCGC
729	21	1	1	1	1	1	GCAGACCGTGCCGGAGGCGCG
730	21	1	1	1	1	1	CAGACCGTGCCGGAGGCGCGC
731	21	1	1	1	1	1	AGACCGTGCCGGAGGCGCGCA
732	21	1	1	1	1	1	GACCGTGCCGGAGGCGCGCAG
733	21	1	1	1	1	1	ACCGTGCCGGAGGCGCGCAGC
734	21	1	1	1	1	1	CCGTGCCGGAGGCGCGCAGCC

59/156

FIG. 20 (36)

[illegible]

60/156

FIG. 20 (37)

757	21	1	2	2	2	2	2	2	GACGCCACGCCCGCGGTGTCC
758	21	1	2	2	2	2	2	2	ACGCCACGCCCGCGGTGTCCC
759	21	1	2	2	2	2	2	2	CGCCACGCCCGCGGTGTCCCC
760	21	1	2	2	2	2	2	2	GCCACGCCCGCGGTGTCCCC
761	21	1	2	2	2	2	2	2	CCACGCCCGCGGTGTCCCCCA
762	21	1	2	2	2	2	2	2	CAGCCGCCCGGTGTCCCCCAT
763	21	1	2	2	2	2	2	2	ACGCCGCCCGGTGTCCCCCATC
764	21	1	2	2	2	2	2	2	CGCCGCCCGGTGTCCCCCATCA
765	21	1	2	2	2	2	2	2	GCCGCCCGGTGTCCCCCATCAA
766	21	1	2	2	2	2	2	2	CCGCCCGGTGTCCCCCATCAAC
767	21	1	2	2	2	2	2	2	CGCCGGTGTCCCCCATCAACA
768	21	1	2	2	2	2	2	2	GCCGGTGTCCCCCATCAACAT
769	21	1	2	2	2	2	2	2	CCGGTGTCCCCCATCAACATG
770	21	1	2	2	2	2	2	2	CGGTGTCCCCCATCAACATGG
771	21	1	2	2	2	2	2	2	GGTGTCCCCCATCAACATGGA
772	21	2	2	2	2	2	2	2	GTGTCCCCCATCAACATGGAA
773	21	2	2	2	2	2	2	2	TGTCCCCCATCAACATGGAG
774	21	2	2	2	2	2	2	2	GTCCCCCATCAACATGGAAGA
775	21	2	2	2	2	2	2	2	TCCCCCATCAACATGGAAGAC
776	21	2	2	2	2	2	2	2	CCCCCATCAACATGGAAGACC
777	21	2	2	2	2	2	2	2	CCCCATCAACATGGAAGACCA

61/156

FIG. 20 (38)

[illegible]



62/156

FIG. 20 (39)

799	21	1	2	2	2	2	GAGCGCATCAAGTGGAGCGC
800	21	1	2	2	2	2	AGCGCATCAAGTGGAGCGCA
801	21	1	2	2	2	2	GCGCATCAAGTGGAGCGCAA
802	21	1	2	2	2	2	CGCATCAAGTGGAGCGCAAG
803	21	1	2	2	2	2	GCATCAAGTGGAGCGCAAGC
804	21	1	2	2	2	2	CATCAAGTGGAGCGCAAGCG
805	21	1	2	2	2	2	ATCAAGTGGAGCGCAAGCGG
806	21	1	2	2	2	2	TCAAAGTGGAGCGCAAGCGGC
807	21	1	2	2	2	2	CAAAGTGGAGCGCAAGCGGCT
808	21	1	2	2	2	2	AAAGTGGAGCGCAAGCGGCTG
809	21	1	2	2	2	2	AAGTGGAGCGCAAGCGGCTGC
810	21	1	2	2	2	2	AGTGGAGCGCAAGCGGCTGCG
811	21	1	2	2	2	2	GTGGAGCGCAAGCGGCTGCGG
812	21	1	2	2	2	2	TGGAGCGCAAGCGGCTGCGGA
813	21	1	2	2	2	2	GGAGCGCAAGCGGCTGCGGAA
814	21	1	2	2	2	2	GAGCGCAAGCGGCTGCGGAAC
815	21	1	1	1	1	1	AGCGCAAGCGGCTGCGGAACC
816	21	1	1	1	1	1	GCGCAAGCGGCTGCGGAACCG
817	21	1	1	1	1	1	CGCAAGCGGCTGCGGAACCGG
818	21	1	1	1	1	1	GCAAGCGGCTGCGGAACCGGC
819	21	1	1	1	1	1	CAAGCGGCTGCGGAACCGGCT
820	21	1	2	2	2	2	AAGCGGCTGCGGAACCGGCTG

63/156

FIG. 20 (40)

821	21	1	2	2	2	2	2	AGGGCTGCGGAACCGGCTGG
822	21	1	2	2	2	2	2	GCGGCTGCGGAACCGGCTGGC
823	21	1	2	2	2	2	2	CGGCTGCGGAACCGGCTGGCG
824	21	1	2	2	2	2	2	GGCTGCGGAACCGGCTGGCGG
825	21	1	2	2	2	2	2	GCTGCGGAACCGGCTGGCGGC
826	21	1	2	2	2	2	2	CTGCGGAACCGGCTGGCGGCC
827	21	1	2	2	2	2	2	TGCGGAACCGGCTGGCGGCCA
828	21	1	2	2	2	2	2	GCGAACCGGCTGGCGGCCAC
829	21	1	2	2	2	2	2	CGBAACCGGCTGGCGGCCACC
830	21	1	2	2	2	2	2	GGAAACCGGCTGGCGGCCACCA
831	21	1	2	2	2	2	2	GAACCGGCTGGCGGCCACCAA
832	21	1	2	2	2	2	2	AACCGGCTGGCGGCCACCAAG
833	21	1	2	2	2	2	2	ACCGGCTGGCGGCCACCAAGT
834	21	1	2	2	2	2	2	CCGGCTGGCGGCCACCAAGTG
835	21	1	2	2	2	2	2	CGGCTGGCGGCCACCAAGTGC
836	21	2	2	2	2	2	2	GGCTGGCGGCCACCAAGTGCC
837	21	2	2	2	2	2	2	GCTGGCGGCCACCAAGTGCCG
838	21	2	2	2	2	2	2	CTGGCGGCCACCAAGTGCCGG
839	21	2	2	2	2	2	2	TGGCGGCCACCAAGTGCCGGA
840	21	2	2	2	2	2	2	GGCGGCCACCAAGTGCCGGAA
841	21	2	2	2	2	2	2	GCGGCCACCAAGTGCCGGGAA

64/156

FIG. 20 (41)

842	21	2	2	2	2	2	CGGCCACCAAGTGCCGGAAGC
843	21	2	2	2	2	2	GGCCACCAAGTGCCGGAAGCG
844	21	2	2	2	2	2	GCCACCAAGTGCCGGAAGCGG
845	21	2	2	2	2	2	CCACCAAGTGCCGGAAGCGGA
846	21	2	2	2	2	2	CACCAAGTGCCGGAAGCGGAA
847	21	2	2	2	2	2	ACCAAGTGCCGGAAGCGGAAG
848	21	2	2	2	2	2	CCAAGTGCCGGAAGCGGAAGC
849	21	2	2	2	2	2	CAAGTGCCGGAAGCGGAAGCT
850	21	2	2	2	2	2	AAGTGCCGGAAGCGGAAGCTG
851	21	2	2	2	2	2	AGTGCCGGAAGCGGAAGCTGG
852	21	2	2	2	2	2	GTGCCGGAAGCGGAAGCTGGA
853	21	2	2	2	2	2	TGCCGGAAGCGGAAGCTGGAG
854	21	2	2	2	2	2	GCCGGAAGCGGAAGCTGGAGC
855	21	2	2	2	2	2	CCGGAAGCGGAAGCTGGAGCG
856	21	2	2	2	2	2	CGGAAGCGGAAGCTGGAGCGC
857	21	2	2	2	2	2	GGAAGCGGAAGCTGGAGCGCA
858	21	2	2	2	2	2	GAAGCGGAAGCTGGAGCGCAT
859	21	2	2	2	2	2	AAGCGGAAGCTGGAGCGCATC
860	21	2	2	2	2	2	AGCGGAAGCTGGAGCGCATCG
861	21	2	2	2	2	2	GCGGAAGCTGGAGCGCATCGC
862	21	2	2	2	2	2	CGGAAGCTGGAGCGCATCGCG

65/156

FIG. 20 (42)

[illegible]

66/156

FIG. 20 (43)

885	21	2	2	2	2	2	CCTGGAGGACAAGGTGAAGAC
886	21	2	2	2	2	2	CTGAGGACAAGGTGAAGACG
887	21	1	2	2	2	2	TGGAGGACAAGGTGAAGACGC
888	21	1	2	2	2	2	GGAGGACAAGGTGAAGACGCT
889	21	1	2	2	2	2	GAGGACAAGGTGAAGACGCTC
890	21	1	2	2	2	2	AGGACAAGGTGAAGACGCTCA
891	21	1	2	2	2	2	GGACAAGGTGAAGACGCTCAA
892	21	1	2	2	2	2	GACAAGGTGAAGACGCTCAAG
893	21	1	2	2	2	2	ACAAGGTGAAGACGCTCAAGG
894	21	1	2	2	2	2	CAAGGTGAAGACGCTCAAGGC
895	21	1	1	2	1	1	AAGGTGAAGACGCTCAAGGCC
896	21	1	1	1	1	1	AGGTGAAGACGCTCAAGGCCG
897	21	1	1	1	1	1	GGTGAAGACGCTCAAGGCCGA
898	21	1	1	1	1	1	GTGAAGACGCTCAAGGCCGAG
899	21	1	1	1	1	1	TGAAGACGCTCAAGGCCGAGA
900	21	1	1	1	1	1	GAAGACGCTCAAGGCCGAGAA
901	21	1	1	1	1	1	AAGACGCTCAAGGCCGAGAAC
902	21	1	1	1	1	1	AGACGCTCAAGGCCGAGAACG
903	21	1	1	1	1	1	GACGCTCAAGGCCGAGAACGC
904	21	1	1	1	1	1	ACGCTCAAGGCCGAGAACGCG
905	21	1	1	1	1	1	CGCTCAAGGCCGAGAACGCGG

67/156

FIG. 20 (44)

906	21	1	1	1	1	1	GCTCAAGGCCGAGAACGCGGG
907	21	1	2	2	2	2	CTCAAGGCCGAGAACGCGGG
908	21	1	2	2	2	2	TCAAGGCCGAGAACGCGGGC
909	21	1	2	2	2	2	CAAGGCCGAGAACGCGGGCT
910	21	1	2	2	2	2	AAGGCCGAGAACGCGGGCTG
911	21	1	2	2	2	2	AGGCCGAGAACGCGGGCTGT
912	21	1	2	2	2	2	GGCCGAGAACGCGGGGCTGTC
913	21	1	2	2	2	2	GCCGAGAACGCGGGGCTGTCG
914	21	1	2	2	2	2	CCGAGAACGCGGGGCTGTCCA
915	21	1	2	2	2	2	CGAGAACGCGGGGCTGTGAG
916	21	2	2	2	2	2	GAGAACGCGGGGCTGTGAGT
917	21	1	2	2	2	2	AGAACGCGGGGCTGTGAGTA
918	21	1	2	2	2	2	GAACGCGGGGCTGTGAGTAC
919	21	1	1	1	1	1	AACGCGGGGCTGTGAGTACC
920	21	1	1	1	1	1	ACGCGGGGCTGTGAGTACCG
921	21	1	1	1	1	1	CGCGGGGCTGTGAGTACCGC
922	21	1	1	1	1	1	GCGGGGCTGTGAGTACCGCC
923	21	1	1	1	1	1	CGGGGCTGTGAGTACCGCCG
924	21	1	1	1	1	1	GGGGCTGTGAGTACCGCCGG
925	21	1	1	1	1	1	GGGCTGTGAGTACCGCCGGC
926	21	1	1	1	1	1	GGCTGTGAGTACCGCCGGCC

68/156

FIG. 20 (45)

927	1	1	1	1	1	GCTGTCGAGTACCGCCGGCCT
928	1	1	1	1	1	CTGTCGAGTACCGCCGGCCTC
929	1	1	1	1	1	TGTCGAGTACCGCCGGCCTCC
930	1	1	1	1	1	GTCGAGTACCGCCGGCCTCCT
931	1	1	1	1	1	TCGAGTACCGCCGGCCTCCTC
932	1	1	1	1	1	CGAGTACCGCCGGCCTCCTCC
933	1	1	1	1	1	GAGTACCGCCGGCCTCCTCCG
934	1	1	1	1	1	AGTACCGCCGGCCTCCTCCGG
935	1	1	1	1	1	GTACCGCCGGCCTCCTCCGGG
936	1	1	1	1	1	TACCGCCGGCCTCCTCCGGGA
937	1	1	1	1	1	ACCGCCGGCCTCCTCCGGGAG
938	1	1	1	1	1	CCGCCGGCCTCCTCCGGGAGC
939	1	1	1	1	1	CGCCGGCCTCCTCCGGGAGCA
940	1	1	1	1	1	GCCGGCCTCCTCCGGGAGCAG
941	1	1	1	1	1	CCGGCCTCCTCCGGGAGCAGG
942	1	1	1	1	1	CGGCCTCCTCCGGGAGCAGGT
943	1	1	1	1	1	GGCCTCCTCCGGGAGCAGGTG
944	1	1	1	1	1	GCCTCCTCCGGGAGCAGGTGG
945	1	1	1	1	1	CCTCCTCCGGGAGCAGGTGGC
946	1	1	1	1	1	CTCCTCCGGGAGCAGGTGGCC
947	1	1	1	1	1	TCCTCCGGGAGCAGGTGGCCC
948	1	1	1	1	1	CCTCCGGGAGCAGGTGGCCCA

69/156

FIG. 20 (46)

949	21	1	1	1	1	1	CTCCGGGAGCAGGTGGCCCAG
950	21	1	1	1	1	1	TCCGGGAGCAGGTGGCCCAGC
951	21	1	1	1	1	1	CCGGAGCAGGTGGCCCAGCT
952	21	1	1	1	1	1	CGGAGCAGGTGGCCCAGCTC
953	21	1	1	1	1	1	GGAGCAGGTGGCCCAGCTCA
954	21	1	1	1	1	1	GGAGCAGGTGGCCCAGCTCAA
955	21	1	2	2	2	2	GAGCAGGTGGCCCAGCTCAA
956	21	1	2	2	2	2	AGCAGGTGGCCCAGCTCAAAC
957	21	1	2	2	2	2	GCAGGTGGCCCAGCTCAAACA
958	21	1	2	2	2	2	CAGGTGGCCCAGCTCAAACAG
959	21	1	2	2	2	2	AGGTGGCCCAGCTCAAACAGA
960	21	1	2	2	2	2	GGTGGCCCAGCTCAAACAGAA
961	21	1	1	1	1	1	GTGGCCCAGCTCAAACAGAAAG
962	21	1	1	1	1	1	TGGCCCAGCTCAAACAGAAAGG
963	21	1	1	1	1	1	GGCCCAGCTCAAACAGAAAGGT
964	21	1	1	1	1	1	GGCCAGCTCAAACAGAAAGGTC
965	21	1	1	1	1	1	CCAGCTCAAACAGAAAGGTCA
966	21	1	1	1	1	1	CCAGCTCAAACAGAAAGGTCA
967	21	1	2	2	2	2	CAGCTCAAACAGAAAGGTCA
968	21	1	2	2	2	2	CAGCTCAAACAGAAAGGTCA
969	21	1	2	2	2	2	AGCTCAAACAGAAAGGTCA
		1	2	2	2	2	GCTCAAACAGAAAGGTCA
							GAC



70/156

FIG. 20 (47)

970	21	1	2	2	2	2	2	2	CTCAACAGAAAGGTCATGACC
971	21	1	2	2	2	2	2	2	TCAACAGAAAGGTCATGACCC
972	21	1	2	2	2	2	2	2	CAACAGAAAGGTCATGACCCA
973	21	1	1	1	1	1	1	1	AAACAGAAAGGTCATGACCCAC
974	21	1	1	1	1	1	1	1	AACAGAAAGGTCATGACCCACG
975	21	1	1	1	1	1	1	1	ACAGAAAGGTCATGACCCACGT
976	21	1	2	2	2	2	2	2	CAGAAGGTCATGACCCACGTC
977	21	1	2	2	2	2	2	2	AGAAGGTCATGACCCACGTCA
978	21	1	2	2	2	2	2	2	GAAGGTCATGACCCACGTGTCAG
979	21	1	2	2	2	2	2	2	AAGGTCATGACCCACGTGTCAGC
980	21	1	2	2	2	2	2	2	AGGTCATGACCCACGTGTCAGCA
981	21	1	2	2	2	2	2	2	GGTCATGACCCACGTGTCAGCAA
982	21	1	2	2	2	2	2	2	GTCATGACCCACGTGTCAGCAAC
983	21	1	2	2	2	2	2	2	TCATGACCCACGTGTCAGCAACG
984	21	1	2	2	2	2	2	2	CATGACCCACGTGTCAGCAACGG
985	21	1	2	2	2	2	2	2	ATGACCCACGTGTCAGCAACGGC
986	21	1	2	2	2	2	2	2	TGACCCACGTGTCAGCAACGGCT
987	21	1	2	2	2	2	2	2	GACCCACGTGTCAGCAACGGCTG
988	21	1	1	1	1	1	1	1	ACCCACGTGTCAGCAACGGCTGT
989	21	1	1	1	1	1	1	1	CCCACGTGTCAGCAACGGCTGTC
990	21	1	1	1	1	1	1	1	CCACGTGTCAGCAACGGCTGTCA



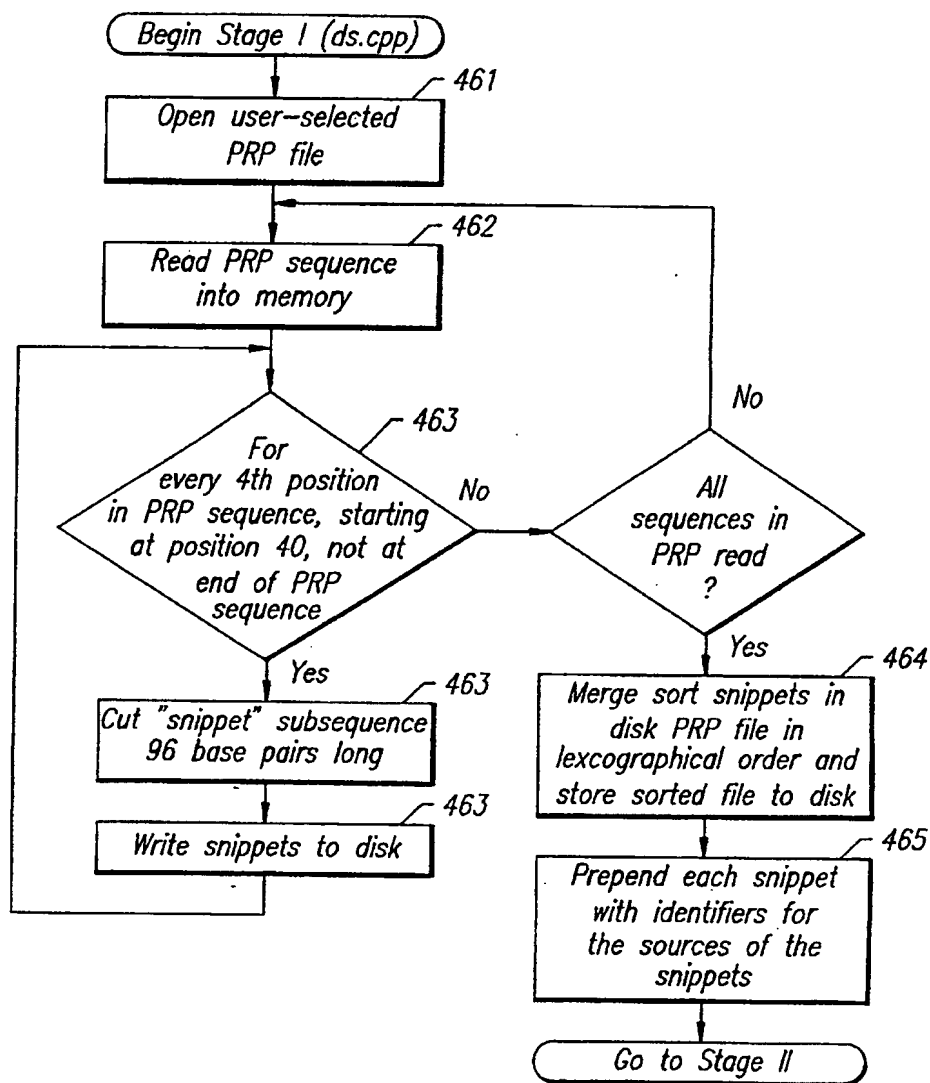
72/156

FIG. 20 (49)

11013	21	1	2	2	2	2	TGCTGCTTGGGGGTCAAGGGAC
11014	21	1	2	2	2	2	GCTGCTTGGGGGTCAAGGGACA
11015	21	1	2	2	2	2	CTGCTTGGGGGTCAAGGGACAC
11016	21	1	2	2	2	2	TGCTTGGGGTCAAGGGACACG
11017	21	1	2	2	2	2	GCTTGGGGTCAAGGGACACGC
11018	21	1	2	2	2	2	CTTGGGGTCAAGGGACACGCC
11019	21	1	2	2	2	2	TTGGGGTCAAGGGACACGCCT
11020	21	1	2	2	2	2	TGGGGTCAAGGGACACGCCCTT
11021	21	2	2	2	2	2	GGGTCAAGGGACACGCCCTTC
11022	21	2	2	2	2	2	GGGTCAAGGGACACGCCCTTCT
11023	21	2	2	2	2	2	GGTCAAGGGACACGCCCTTCTG
11024	21	2	2	2	2	2	GTCAAGGGACACGCCCTTCTGA

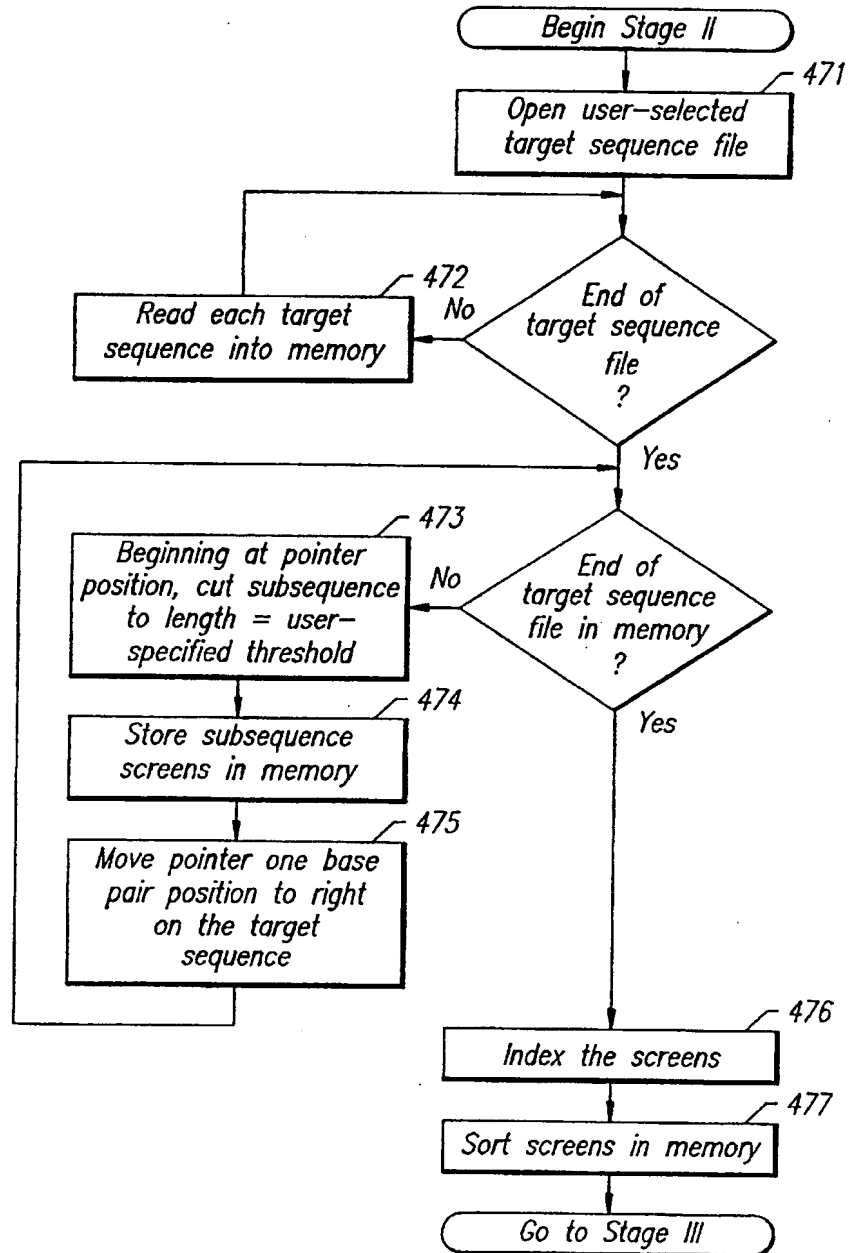
73/156

FIG. 21



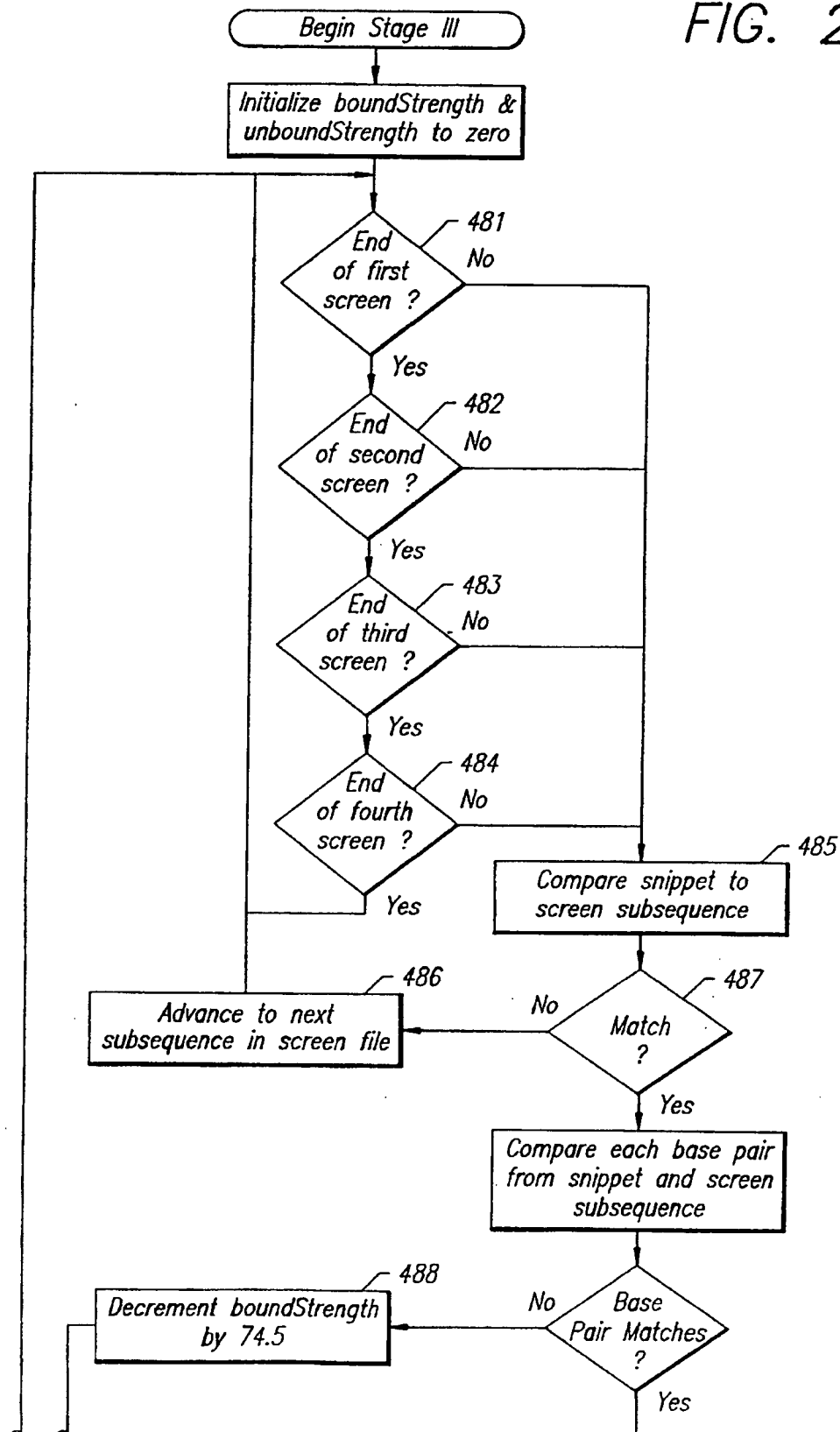
74/156

FIG. 22



75/156

FIG. 23(1)



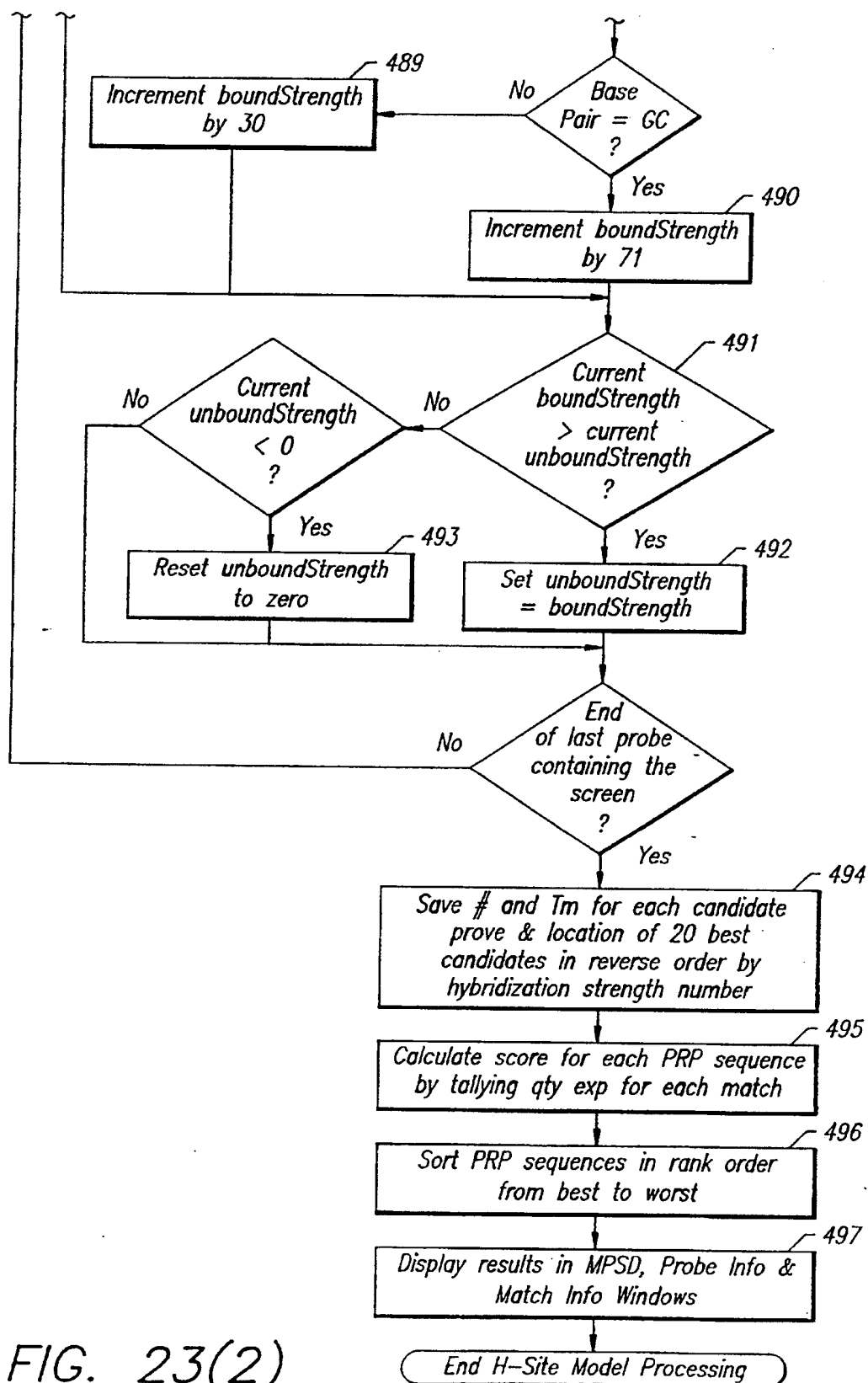


FIG. 23(2)

## FIG. 24A (1)

## OligoProbe DesignStation

Probes: C:\HITACHI\HUMBJUNX.CDS  
 Datatbase: C:\HITACHI\JUNMIX.SEQ

Mismatch Model, 1 = 21, k = 4

Position length	Mismatches								screensN Probe	77/156
	0	1	2	3	4	5	6	7	8	
1 21	0	0	0	0	0	0	0	0	0	ATGTGCACTAAAAATGGAACAG
2 21	0	0	0	0	0	0	0	0	0	TGTGCACTAAAAATGGAACAGC
3 21	0	0	0	0	0	0	0	0	0	GTGCACTAAAAATGGAACAGCC
4 21	0	0	0	0	0	0	0	0	0	TGCACATAAAATGGAACAGCCC
5 21	0	0	0	0	0	0	0	0	0	GCACTAAAAATGGAACAGCCCT
6 21	0	0	0	0	0	0	0	0	0	CACTAAAAATGGAACAGCCCTT
7 21	0	0	0	0	0	0	0	0	0	ACTAAAAATGGAACAGCCCTTC
8 21	0	0	0	0	0	0	0	0	0	CTAAAAATGGAACAGCCCTTCT
9 21	0	0	0	0	0	0	0	0	0	TAAAAATGGAACAGCCCTTCTA
10 21	0	0	0	0	0	0	0	0	0	AAAAATGGAACAGCCCTTCTAC



78/156

FIG. 24A (2)

11	21	0	0	0	0	0	0	AAATGGAACAGCCCTTCTACC
12	21	0	0	0	0	0	0	AATGGAACAGCCCTTCTACCA
13	21	0	0	0	0	0	0	ATGGAACAGCCCTTCTACCA
14	21	0	0	0	0	0	0	TGGAACAGCCCTTCTACCA
15	21	0	0	0	0	0	0	GGAACAGCCCTTCTACCA
16	21	0	0	0	0	0	0	GAAACAGCCCTTCTACCA
17	21	0	0	0	0	0	0	AACAGCCCTTCTACCA
18	21	0	0	0	0	0	0	ACAGCCCTTCTACCA
19	21	0	0	0	0	0	0	CAGCCCTTCTACCA
20	21	0	0	0	0	0	0	AGCCCTTCTACCA
21	21	0	0	0	0	0	0	GCCCTTCTACCA
22	21	0	0	0	0	0	0	CCCTTCTACCA
23	21	0	0	0	0	0	0	CCTTCTACCA
24	21	0	0	0	0	0	0	CTTCTACCA
25	21	0	0	0	0	0	0	TTCTACCA
26	21	0	0	0	0	0	0	TCTACCA
27	21	0	0	0	0	0	0	CTACCA
28	21	0	0	0	0	0	0	TACCA
29	21	0	0	0	0	0	0	ACCAC
30	21	0	0	0	0	0	0	CCAC
31	21	0	0	0	0	0	0	CAC









FIG. 24A (7)

[illegible]

84/156

FIG. 24A. (8)

[illegible]

85/156

FIG. 24A (9)

158	21	0	0	0	0	0	0	0	CTGGGGCTCGCGGACCCGGGCC
159	21	0	0	0	0	0	0	0	TGGGGCTCGCGGACCCGGGCC
160	21	0	0	0	0	0	0	0	GGGGCTCGCGGACCCGGGCCA
161	21	0	0	0	0	0	0	0	GGGCTCGCGGACCCGGGCCAG
162	21	0	0	0	0	0	0	0	GGCTCGCGGACCCGGGCCAGA
163	21	0	0	0	0	0	0	0	GCTCGCGGACCCGGGCCAGAG
164	21	0	0	0	0	0	0	0	CTCGGGACCCGGGCCAGAGG
165	21	0	0	0	0	0	0	0	TCGCGGACCCGGGCCAGAGGG
166	21	0	0	0	0	0	0	0	CGCGACCCGGGCCAGAGGGC
167	21	0	0	0	0	0	0	0	GCGGACCCGGGCCAGAGGGCG
168	21	0	0	0	0	0	0	0	CGGACCCGGGCCAGAGGGCGG
169	21	0	0	0	0	0	0	0	GGACCCGGGCCAGAGGGCGGC
170	21	0	0	0	0	0	0	0	GACCCGGGCCAGAGGGCGGCG
171	21	0	0	0	0	0	0	0	ACCCGGGCCAGAGGGCGGCGG
172	21	0	0	0	0	0	0	0	CCGGCCAGAGGGCGGCGGCT
173	21	0	0	0	0	0	0	0	CCGGCCAGAGGGCGGCGGTG
174	21	0	0	0	0	0	0	0	CGGCCAGAGGGCGGCGGTGG
175	21	0	0	0	0	0	0	0	GGCCAGAGGGCGGCGGTGGC
176	21	0	0	0	0	0	0	0	GCCAGAGGGCGGCGGTGGCG
177	21	0	0	0	0	0	0	0	CCAGAGGGCGGCGGTGGCGG
178	21	0	0	0	0	0	0	0	CCAGAGGGCGGCGGTGGCGGC



86/156

FIG. 24A (10)

[illegible]

87/156

FIG. 24A (11)

200	21	0	0	0	0	0	GCTACTTTTCTGTCAGGGCT
201	21	0	0	0	0	0	CTACTTTTCTGTCAGGGCTC
202	21	0	0	0	0	0	TACTTTTCTGTCAGGGCTCG
203	21	0	0	0	0	0	ACTTTTCTGTCAGGGCTCGG
204	21	0	0	0	0	0	CTTTTCTGTCAGGGCTCGGA
205	21	0	0	0	0	0	TTTTTCTGTCAGGGCTCGGAC
206	21	0	0	0	0	0	TTTCTGGTCAGGGCTCGGACA
207	21	0	0	0	0	0	TTCTGGTCAGGGCTCGGACAC
208	21	0	0	0	0	0	TCTGGTCAGGGCTCGGACACC
209	21	0	0	0	0	0.	CTGGTCAGGGCTCGGACACCG
210	21	0	0	0	0	0	TGGTCAGGGCTCGGACACCCGG
211	21	0	0	0	0	0	GGTCAGGGCTCGGACACCCGGC
212	21	0	0	0	0	0	GTCAGGGCTCGGACACCCGGCG
213	21	0	0	0	0	0	TCAGGGCTCGGACACCCGGCGC
214	21	0	0	0	0	0	CAGGGCTCGGACACCCGGCGCG
215	21	0	0	0	0	0	AGGGCTCGGACACCCGGCGCGT
216	21	0	0	0	0	0	GGGCTCGGACACCCGGCGCGTC
217	21	0	0	0	0	0	GGCTCGGACACCCGGCGGTCT
218	21	0	0	0	0	0	GCTCGGACACCCGGCGGTCTC
219	21	0	0	0	0	0	CTCGGACACCCGGCGGTCTCT
220	21	0	0	0	0	0	TCGGACACCCGGCGGTCTCTC

88/156

FIG. 24A (12)

221	21	0	0	0	0	0	0	0	CGGACACCGGCGGTCTCTCA
222	21	0	0	0	0	0	0	0	GGACACCGGCGGTCTCTCAA
223	21	0	0	0	0	0	0	0	GACACCGGCGGTCTCTCAAG
224	21	0	0	0	0	0	0	0	ACACCGGCGGTCTCTCAAGC
225	21	0	0	0	0	0	0	0	CACCGGCGGTCTCTCAAGCT
226	21	0	0	0	0	0	0	0	ACCGGCGGTCTCTCAAGCTC
227	21	0	0	0	0	0	0	0	CCGGCGGTCTCTCAAGCTCG
228	21	0	0	0	0	0	0	0	CGGCGGTCTCTCAAGCTCGC
229	21	0	0	0	0	0	0	0	GGCGGTCTCTCAAGCTCGCC
230	21	0	0	0	0	0	0	0	GCGGTCTCTCAAGCTCGCCT
231	21	0	0	0	0	0	0	0	CGGTCTCTCAAGCTCGCCTC
232	21	0	0	0	0	0	0	0	GCGTCTCTCAAGCTCGCCTCT
233	21	0	0	0	0	0	0	0	CGTCTCTCAAGCTCGCCTCTT
234	21	0	0	0	0	0	0	0	GTCCTCTCAAGCTCGCCTCTTC
235	21	0	0	0	0	0	0	0	TCTCTCAAGCTCGCCTCTTCG
236	21	0	0	0	0	0	0	0	CTCTCAAGCTCGCCTCTTCGG
237	21	0	0	0	0	0	0	0	TCTCAAGCTCGCCTCTTCGGA
238	21	0	0	0	0	0	0	0	CTCAAGCTCGCCTCTTCGGAG
239	21	0	0	0	0	0	0	0	TCAAGCTCGCCTCTTCGGAGC
240	21	0	0	0	0	0	0	0	CAAGCTCGCCTCTTCGGAGCT
241	21	0	0	0	0	0	0	0	AAGCTCGCCTCTTCGGAGCTG

89/156

FIG. 24A (13)

242	21	0	0	0	0	0	AGCTCGCCCTCTTCGGAGCTGG
243	21	0	0	0	0	0	GCTCGCCTCTTTCGGAGCTGGA
244	21	0	0	0	0	0	CTCGCCTCTTTCGGAGCTGGAA
245	21	0	0	0	0	0	TGCCTCTTTCGGAGCTGGAAC
246	21	0	0	0	0	0	CGCCCTCTTCGGAGCTGGAACG
247	21	0	0	0	0	0	GCCTCTTCGGAGCTGGAACGC
248	21	0	0	0	0	0	CCTCTTCGGAGCTGGAACGCC
249	21	0	0	0	0	0	CTCTTCGGAGCTGGAACGCCCT
250	21	0	0	0	0	0	TCTTCGGAGCTGGAACGCCCTG
251	21	0	0	0	0	0	CTTCGGAGCTGGAACGCCCTGA
252	21	0	0	0	0	0	TTTCGGAGCTGGAACGCCCTGAT
253	21	0	0	0	0	0	TCGGAGCTGGAACGCCCTGATT
254	21	0	0	0	0	0	CGGAGCTGGAACGCCCTGATTG
255	21	0	0	0	0	0	GGAGCTGGAACGCCCTGATTGT
256	21	0	0	0	0	0	GAGCTGGAACGCCCTGATTGTC
257	21	0	0	0	0	0	AGCTGGAACGCCCTGATTGTCC
258	21	0	0	0	0	0	GCTGGAACGCCCTGATTGTCCC
259	21	0	0	0	0	0	CTGGAACGCCCTGATTGTCCCC
260	21	0	0	0	0	0	TGGAACGCCCTGATTGTCCCCA
261	21	0	0	0	0	0	GGAACGCCCTGATTGTCCCCAA
262	21	0	0	0	0	0	GAACGCCCTGATTGTCCCCCAAC

90/156

FIG. 24A (14)

263	21	0	0	0	0	0	AACGCTGATTGTCCCCAACA
264	21	0	0	0	0	0	ACGCTGATTGTCCCCAACAG
265	21	0	0	0	0	0	CGCTGATTGTCCCCAACAGC
266	21	0	0	0	0	0	GCTGATTGTCCCCAACAGCA
267	21	0	0	0	0	0	CCTGATTGTCCCCAACAGCAA
268	21	0	0	0	0	0	CTGATTGTCCCCAACAGCAAC
269	21	0	0	0	0	0	TGATTGTCCCCAACAGCAACG
270	21	0	0	0	0	0	GATTGTCCCCAACAGCAACGG
271	21	0	0	0	0	0	ATTGTCCCCAACAGCAACGGC
272	21	0	0	0	0	0	TTGTCCCCAACAGCAACGGCG
273	21	0	0	0	0	0	TGTCCCCAACAGCAACGGCGT
274	21	0	0	0	0	0	GTCCCCAACAGCAACGGCGTG
275	21	0	0	0	0	0	TCCCCAACAGCAACGGCGTGA
276	21	0	0	0	0	0	CCCCAACAGCAACGGCGTGAT
277	21	0	0	0	0	0	CCCAACAGCAACGGCGTGATC
278	21	0	0	0	0	0	CCAACAGCAACGGCGTGATCA
279	21	0	0	0	0	0	CAACAGCAACGGCGTGATCAC
280	21	0	0	0	0	0	AACAGCAACGGCGTGATCACG
281	21	0	0	0	0	0	ACAGCAACGGCGTGATCACGA
282	21	0	0	0	0	0	CAGCAACGGCGTGATCACGAC
283	21	0	0	0	0	0	AGCAACGGCGTGATCACGACC

91/156

FIG. 24A (15)

284	21	0	0	0	0	0	GCAACGGCGTGATCACGACGA
285	21	0	0	0	0	0	CAACGGCGTGATCACGACGAC
286	21	0	0	0	0	0	AACGGCGTGATCACGACGACG
287	21	0	0	0	0	0	ACGGCGTGATCACGACGACGC
288	21	0	0	0	0	0	CGGCGTGATCACGACGACGCC
289	21	0	0	0	0	0	GGCGTGATCACGACGACGCCT
290	21	0	0	0	0	0	GCGTGATCACGACGACGCCTA
291	21	0	0	0	0	0	CGTGATCACGACGACGCCTAC
292	21	0	0	0	0	0	GTGATCACGACGACGCCTACA
293	21	0	0	0	0	0	TGATCACGACGACGCCTACAC
294	21	0	0	0	0	0	GATCACGACGACGCCTACACC
295	21	0	0	0	0	0	ATCACGACGACGCCTACACCC
296	21	0	0	0	0	0	TCACGACGACGCCTACACCCC
297	21	0	0	0	0	0	CACGACGACGCCTACACCCCC
298	21	0	0	0	0	0	ACGACGACGCCTACACCCCCG
299	21	0	0	0	0	0	CGACGACGCCTACACCCCCCG
300	21	0	0	0	0	0	GACGACGCCTACACCCCCCGG
301	21	0	0	0	0	0	ACGACGCCTACACCCCCCGGG
302	21	0	0	0	0	0	CGACGCCTACACCCCCCGGGA
303	21	0	0	0	0	0	GACGCCTACACCCCCCGGGAC
304	21	0	0	0	0	0	ACGCCTACACCCCCCGGGACAG

92/156

FIG. 24A (16)

305	21	0	0	0	0	0	0	0	CGCTACACCCCGGACAGT
306	21	0	0	0	0	0	0	0	GCCTACACCCCGGACAGTA
307	21	0	0	0	0	0	0	0	CCTACACCCCGGACAGTAC
308	21	0	0	0	0	0	0	0	CTACACCCCGGACAGTACT
309	21	0	0	0	0	0	0	0	TACACCCCGGACAGTACTT
310	21	0	0	0	0	0	0	0	ACACCCCGGACAGTACTTT
311	21	0	0	0	0	0	0	0	CACCCCGGACAGTACTTTT
312	21	0	0	0	0	0	0	0	ACCCCGGACAGTACTTTTA
313	21	0	0	0	0	0	0	0	CCCCGGACAGTACTTTTAC
314	21	0	0	0	0	0	0	0	CCCCGGACAGTACTTTTACC
315	21	0	0	0	0	0	0	0	CCCGGACAGTACTTTTACCC
316	21	0	0	0	0	0	0	0	CCGGACAGTACTTTTACCCC
317	21	0	0	0	0	0	0	0	CGGACAGTACTTTTACCCCC
318	21	0	0	0	0	0	0	0	GGACAGTACTTTTACCCCCG
319	21	0	0	0	0	0	0	0	GGACAGTACTTTTACCCCCG
320	21	0	0	0	0	0	0	0	GACAGTACTTTTACCCCCCG
321	21	0	0	0	0	0	0	0	ACAGTACTTTTACCCCCCGG
322	21	0	0	0	0	0	0	0	CAGTACTTTTACCCCCCGGG
323	21	0	0	0	0	0	0	0	AGTACTTTTACCCCCCGGGG
324	21	0	0	0	0	0	0	0	GTACTTTTACCCCCCGGGGG
325	21	0	0	0	0	0	0	0	TACTTTTACCCCCCGGGGGT

FIG. 24A (17)

326	21	0	0	0	0	0	0	ACTTTTACCCCGGGGCTG
327	21	0	0	0	0	0	0	CTTTTACCCCGGGGTGG
328	21	0	0	0	0	0	0	TTTTTACCCCGGGGTGCC
329	21	0	0	0	0	0	0	TTTACCCCGGGGGTGCCA
330	21	0	0	0	0	0	0	TTACCCCGGGGGTGGCAG
331	21	0	0	0	0	0	0	TACCCCGGGGGTGGCAGC
332	21	0	0	0	0	0	0	ACCCCGGGGGTGGCAGCG
333	21	0	0	0	0	0	0	CCCCGGGGGTGGCAGCGG
334	21	0	0	0	0	0	0	CCCCGGGGGTGGCAGCGGT
335	21	0	0	0	0	0	0	CCCGGGGGTGGCAGCGGTG
336	21	0	0	0	0	0	0	CCGCGGGGTGGCAGCGGTGG
337	21	0	0	0	0	0	0	CGCGGGGTGGCAGCGGTGGA
338	21	0	0	0	0	0	0	GCGGGGTGGCAGCGGTGCAG
339	21	0	0	0	0	0	0	CGGGGTGGCAGCGGTGGAGG
340	21	0	0	0	0	0	0	GGGGTGGCAGCGGTGGAGGT
341	21	0	0	0	0	0	0	GGGTGGCAGCGGTGGAGGTG
342	21	0	0	0	0	0	0	GGTGGCAGCGGTGGAGGTGC
343	21	0	0	0	0	0	0	GGTGGCAGCGGTGGAGGTGCA
344	21	0	0	0	0	0	0	GTGGCAGCGGTGGAGGTGCAG
345	21	0	0	0	0	0	0	TGGCAGCGGTGGAGGTGCAGG
346	21	0	0	0	0	0	0	GGCAGCGGTGGAGGTGCAGGG



94/156

FIG. 24A (18)

347	21	0	0	0	0	0	0	0	GCAGCGGTGGAGGTGCAGGGG
348	21	0	0	0	0	0	0	0	CAGCGGTGGAGGTGCAGGGGG
349	21	0	0	0	0	0	0	0	AGCGGTGGAGGTGCAGGGGGC
350	21	0	0	0	0	0	0	0	GCGTGGAGGTGCAGGGGGCG
351	21	0	0	0	0	0	0	0	CGGTGGAGGTGCAGGGGGCGC
352	21	0	0	0	0	0	0	0	GGTGGAGGTGCAGGGGGCGCA
353	21	0	0	0	0	0	0	0	GTGAGGTGCAGGGGGCGCAG
354	21	0	0	0	0	0	0	0	TGGAGGTGCAGGGGGCGCAGG
355	21	0	0	0	0	0	0	0	GGAGGTGCAGGGGGCGCAGGG
356	21	0	0	0	0	0	0	0	GAGGTGCAGGGGGCGCAGGGG
357	21	0	0	0	0	0	0	0	AGGTGCAGGGGGCGCAGGGGG
358	21	0	0	0	0	0	0	0	GGTGCAGGGGGCGCAGGGGGC
359	21	0	0	0	0	0	0	0	GTGAGGGGGCGCAGGGGGCG
360	21	0	0	0	0	0	0	0	TGCAGGGGGCGCAGGGGGCGG
361	21	0	0	0	0	0	0	0	GCAGGGGGCGCAGGGGGCGGC
362	21	0	0	0	0	0	0	0	CAGGGGGCGCAGGGGGCGGGC
363	21	0	0	0	0	0	0	0	AGGGGGCGCAGGGGGCGGCGT
364	21	0	0	0	0	0	0	0	GGGGCGCAGGGGGCGGCGTC
365	21	0	0	0	0	0	0	0	GGGGCGCAGGGGGCGGCGTCA
366	21	0	0	0	0	0	0	0	GGGGCGCAGGGGGCGGCGTCAC
367	21	0	0	0	0	0	0	0	GGGGCGCAGGGGGCGGCGTCACC

95/156

FIG. 24A (19)

[illegible]



97/156

FIG. 24A (21)

410	21	0	0	0	0	0	ACGGCTTTGTCAAAAGCCCTGG
411	21	0	0	0	0	0	CGGCTTTGTCAAAGCCCCTGGA
412	21	0	0	0	0	0	GGCTTTGTCAAAGCCCCTGGAC
413	21	0	0	0	0	0	GCTTGTCAAAGGCCCTGGACG
414	21	0	0	0	0	0	CTTTGTCAAAGCCCCCTGGACGA
415	21	0	0	0	0	0	TTTGTCAAAGCCCCTGGACGAT
416	21	0	0	0	0	0	TTGTCAAAGCCCCCTGGACGATC
417	21	0	0	0	0	0	TGTCAAAGCCCCCTGGACGATCT
418	21	0	0	0	0	0	GTCAAAGCCCCCTGGACGATCTG
419	21	0	0	0	0	0	TCAAAGCCCCCTGGACGATCTGC
420	21	0	0	0	0	0	CAAAGCCCCTGGACGATCTGCA
421	21	0	0	0	0	0	AAAGCCCCTGGACGATCTGCAC
422	21	0	0	0	0	0	AAGCCCCTGGACGATCTGCACA
423	21	0	0	0	0	0	AGCCCTGGACGATCTGCACAA
424	21	0	0	0	0	0	GCCCTGGACGATCTGCACAAG
425	21	0	0	0	0	0	CCCTGGACGATCTGCACAAGA
426	21	0	0	0	0	0	CCTGGACGATCTGCACAAGAT
427	21	0	0	0	0	0	CTGGACGATCTGCACAAGATG
428	21	0	0	0	0	0	TGGACGATCTGCACAAGATGA
429	21	0	0	0	0	0	GGACGATCTGCACAAGATGAA
430	21	0	0	0	0	0	GACGATCTGCACAAGATGAAC

98/156

FIG. 24A (22)

[illegible]

99/156

FIG. 24A (23)

[illegible]

100/156

FIG. 24A (24)

473	21	0	0	0	0	0	0	0	CCCTGGGGCGCTACCGGGGGGC
474	21	0	0	0	0	0	0	0	CCTGGGCGCTACCGGGGGGGCC
475	21	0	0	0	0	0	0	0	CTGGGCGCTACCGGGGGGGCCC
476	21	0	0	0	0	0	0	0	TGGGCGCTACCGGGGGGGCCCC
477	21	0	0	0	0	0	0	0	GGGCGCTACCGGGGGGGCCCCC
478	21	0	0	0	0	0	0	0	GGGCTACCGGGGGGGCCCCCG
479	21	0	0	0	0	0	0	0	GCGCTACCGGGGGGGCCCCCGG
480	21	0	0	0	0	0	0	0	CGCTACCGGGGGGGCCCCCGGC
481	21	0	0	0	0	0	0	0	GCTACCGGGGGGGCCCCCGGCT
482	21	0	0	0	0	0	0	0	CTACCGGGGGGGCCCCCGGCTG
483	21	0	0	0	0	0	0	0	TACCGGGGGGGCCCCCGGCTGG
484	21	0	0	0	0	0	0	0	ACCGGGGGGGCCCCCGGCTGGG
485	21	0	0	0	0	0	0	0	CCGGGGGGGGCCCCCGGCTGGGC
486	21	0	0	0	0	0	0	0	CGGGGGGGGGCCCCCGGCTGGGCC
487	21	0	0	0	0	0	0	0	GGGGGGGGGGCCCCCGGCTGGGCC
488	21	0	0	0	0	0	0	0	GGGGGGGGGGCCCCCGGCTGGGCCCG
489	21	0	0	0	0	0	0	0	GGGGGGGGGGCCCCCGGCTGGGCCCGG
490	21	0	0	0	0	0	0	0	GGGGGGGGGGCCCCCGGCTGGGCCCGGG
491	21	0	0	0	0	0	0	0	GGGGGGGGGGCCCCCGGCTGGGCCCGGGG
492	21	0	0	0	0	0	0	0	GCGGGGGGGGGCCCCCGGCTGGGGGG
493	21	0	0	0	0	0	0	0	CCCCGGGGGGGGCCCCCGGCGGGGC

101/156

FIG. 24A (25)

494	21	0	0	0	0	0	0	0	CCCCGGCTGGCCCCCGGGGCGG
495	21	0	0	0	0	0	0	0	CCCGGCTGGGCCCGGGGCGGT
496	21	0	0	0	0	0	0	0	CCGGCTGGGCCCGGGGCGTC
497	21	0	0	0	0	0	0	0	CGGCTGGGCCCGGGGCGTCT
498	21	0	0	0	0	0	0	0	GGCTGGGCCCGGGGCGTCTA
499	21	0	0	0	0	0	0	0	GCTGGGCCCGGGGCGTCTAC
500	21	0	0	0	0	0	0	0	CTGGGCCCGGGGCGTCTACG
501	21	0	0	0	0	0	0	0	TGGGCCCGGGGCGTCTACGC
502	21	0	0	0	0	0	0	0	GGGCCCGGGGCGTCTACGCC
503	21	0	0	0	0	0	0	0	GGCCCGGGGCGTCTACGCCG
504	21	0	0	0	0	0	0	0	GCCCGGGGCGTCTACGCCGG
505	21	0	0	0	0	0	0	0	CCCGGGGCGTCTACGCCGGC
506	21	0	0	0	0	0	0	0	CCGGGGCGTCTACGCCGGCC
507	21	0	0	0	0	0	0	0	CGGGGGCGTCTACGCCGGCC
508	21	0	0	0	0	0	0	0	CGGGCGTCTACGCCGGCCCG
509	21	0	0	0	0	0	0	0	GGGGCGTCTACGCCGGCCCGG
510	21	0	0	0	0	0	0	0	GGGCGTCTACGCCGGCCCGGA
511	21	0	0	0	0	0	0	0	GGCGTCTACGCCGGCCCGGAG
512	21	0	0	0	0	0	0	0	GCGTCTACGCCGGCCCGGAGC
513	21	0	0	0	0	0	0	0	CGTCTACGCCGGCCCGGAGCC
514	21	0	0	0	0	0	0	0	GTCTACGCCGGCCCGGAGCCA





103/156

FIG. 24A (27)

536	21	0	0	0	0	0	CTCCCGTTTACACCAACCTCA
537	21	0	0	0	0	0	TCCCGTTTACACCAACCTCAG
538	21	0	0	0	0	0	CCCGTTTACACCAACCTCAGC
539	21	0	0	0	0	0	CCGTTTACACCAACCTCAGCA
540	21	0	0	0	0	0	CGTTTACACCAACCTCAGCAG
541	21	0	0	0	0	0	GTTTACACCAACCTCAGCAGC
542	21	0	0	0	0	0	TTTACACCAACCTCAGCAGCT
543	21	0	0	0	0	0	TTACACCAACCTCAGCAGCTA
544	21	0	0	0	0	0	TACACCAACCTCAGCAGCTAC
545	21	0	0	0	0	0	ACACCAACCTCAGCAGCTACT
546	21	0	0	0	0	0	CACCAACCTCAGCAGCTACTC
547	21	0	0	0	0	0	ACCAACCTCAGCAGCTACTCC
548	21	0	0	0	0	0	CCAACCTCAGCAGCTACTCCC
549	21	0	0	0	0	0	CAACCTCAGCAGCTACTCCCC
550	21	0	0	0	0	0	AACCTCAGCAGCTACTCCCCA
551	21	0	0	0	0	0	ACCTCAGCAGCTACTCCCCAG
552	21	0	0	0	0	0	CCTCAGCAGCTACTCCCCCAGC
553	21	0	0	0	0	0	CTCAGCAGCTACTCCCCCAGCC
554	21	0	0	0	0	0	TCAGCAGCTACTCCCCCAGCCT
555	21	0	0	0	0	0	CAGCAGCTACTCCCCCAGCCTC
556	21	0	0	0	0	0	AGCAGCTACTCCCCCAGCCTCT

104/156

FIG. 24A (28)

557	21	0	0	0	0	0	GCAGCTACTCCCCAGCCTCTG	0
558	21	0	0	0	0	0	CAGCTACTCCCCAGCCTCTGC	0
559	21	0	0	0	0	0	AGCTACTCCCCAGCCTCTGCG	0
560	21	0	0	0	0	0	GCTACTCCCCAGCCTCTGCGT	0
561	21	0	0	0	0	0	CTACTCCCCAGCCTCTGCGTC	0
562	21	0	0	0	0	0	TACTCCCCAGCCTCTGCGTCC	0
563	21	0	0	0	0	0	ACTCCCCAGCCTCTGCGTCCT	0
564	21	0	0	0	0	0	CTCCCCAGCCTCTGCGTCCTC	0
565	21	0	0	0	0	0	TCCCCAGCCTCTGCGTCCTCG	0
566	21	0	0	0	0	0	CCCCAGCCTCTGCGTCCTCGG	0
567	21	0	0	0	0	0	CCAGCCTCTGCGTCCTCGGG	0
568	21	0	0	0	0	0	CCAGCCTCTGCGTCCTCGGGA	0
569	21	0	0	0	0	0	CAGCCTCTGCGTCCTCGGGAG	0
570	21	0	0	0	0	0	AGCCTCTGCGTCCTCGGGAGG	0
571	21	0	0	0	0	0	GCCTCTGCGTCCTCGGGAGGC	0
572	21	0	0	0	0	0	CCTCTGCGTCCTCGGGAGGCG	0
573	21	0	0	0	0	0	CTCTGCGTCCTCGGGAGGCGC	0
574	21	0	0	0	0	0	TCTGCGTCCTCGGGAGGCGCC	0
575	21	0	0	0	0	0	CTGCGTCCTCGGGAGGCGCCG	0
576	21	0	0	0	0	0	TGCGTCCTCGGGAGGCGCCGG	0
577	21	0	0	0	0	0	CGGTCTCTCGGGAGGCGCCGGG	0

185/156

FIG. 24A (29)

578	21	0	0	0	0	0	0	0	CGTCCTCGGGAGGCGCCGGGG
579	21	0	0	0	0	0	0	0	GTCCTCGGGAGGCGCCGGGGC
580	21	0	0	0	0	0	0	0	TCCTCGGGAGGCGCCGGGGCT
581	21	0	0	0	0	0	0	0	CCTCGGGAGGCGCCGGGGCTG
582	21	0	0	0	0	0	0	0	CTCGGGAGGCGCCGGGGCTGC
583	21	0	0	0	0	0	0	0	TCGGAGGCGCCGGGGCTGCC
584	21	0	0	0	0	0	0	0	CGGAGGCGCCGGGGCTGCCG
585	21	0	0	0	0	0	0	0	GGAGGCGCCGGGGCTGCCGT
586	21	0	0	0	0	0	0	0	GGAGGCGCCGGGGCTGCCGTC
587	21	0	0	0	0	0	0	0	GAGCGCCGGGGCTGCCGTCG
588	21	0	0	0	0	0	0	0	AGCGCCGGGGCTGCCGTCGG
589	21	0	0	0	0	0	0	0	GGCGCCGGGGCTGCCGTCGGG
590	21	0	0	0	0	0	0	0	GGCGGGGGCTGCCGTCGGGA
591	21	0	0	0	0	0	0	0	CGCCGGGGCTGCCGTCGGGAC
592	21	0	0	0	0	0	0	0	GCCGGGGCTGCCGTCGGGACC
593	21	0	0	0	0	0	0	0	CCGGGGCTGCCGTCGGGACCG
594	21	0	0	0	0	0	0	0	CGGGCTGCCGTCGGGACCGG
595	21	0	0	0	0	0	0	0	GGGGCTGCCGTCGGGACCGGG
596	21	0	0	0	0	0	0	0	GGGTGCCGTCGGGACCGGGA
597	21	0	0	0	0	0	0	0	GGTGCCGTCGGGACCGGGAG
598	21	0	0	0	0	0	0	0	GCTGCCGTCGGGACCGGGAGC

106/156

FIG. 24A (30)

599	21	0	0	0	0	0	CTGCCGT
600	21	0	0	0	0	0	CGGACCCGGAGCTC
601	21	0	0	0	0	0	GCCGTCGGGACCCGGAGCTCG
602	21	0	0	0	0	0	CCGTCGGGACCCGGAGCTCGT
603	21	0	0	0	0	0	CGTCGGGACCCGGAGCTCGTA
604	21	0	0	0	0	0	GTCGGGACCCGGAGCTCGTAC
605	21	0	0	0	0	0	TCCGGACCCGGAGCTCGTACC
606	21	0	0	0	0	0	CCGGACCCGGAGCTCGTACCC
607	21	0	0	0	0	0	GGACCCGGAGCTCGTACCCCG
608	21	0	0	0	0	0	GGACCCGGAGCTCGTACCCGA
609	21	0	0	0	0	0	GACCCGGAGCTCGTACCCGAC
610	21	0	0	0	0	0	ACCCGGAGCTCGTACCCGACG
611	21	0	0	0	0	0	CCGGAGCTCGTACCCGACGA
612	21	0	0	0	0	0	CCGGAGCTCGTACCCGACGAC
613	21	0	0	0	0	0	GGAGCTCGTACCCGACGACC
614	21	0	0	0	0	0	GGAGCTCGTACCCGACGACCA
615	21	0	0	0	0	0	GAGCTCGTACCCGACGACCCAC
616	21	0	0	0	0	0	AGCTCGTACCCGACGACCCACC
617	21	0	0	0	0	0	GCTCGTACCCGACGACCCACCA
618	21	0	0	0	0	0	CTCGTACCCGACGACCCACCAT
619	21	0	0	0	0	0	TCGTACCCGACGACCCACCATC





109/156

FIG. 24A (33)

[illegible]



116/156

FIG. 24A (34)

683	21	0	0	0	0	0	0	0	CGCAGCTGGGCTTGGGCGCGG
684	21	0	0	0	0	0	0	0	GCAGCTGGGCTTGGGCGCGG
685	21	0	0	0	0	0	0	0	CAGCTGGGCTTGGGCGCGG
686	21	0	0	0	0	0	0	0	AGCTGGGCTTGGGCGCGG
687	21	0	0	0	0	0	0	0	GCTGGGCTTGGGCGCGG
688	21	0	0	0	0	0	0	0	CTGGGCTTGGGCGCGG
689	21	0	0	0	0	0	0	0	TGGGCTTGGGCGCGG
690	21	0	0	0	0	0	0	0	GGGCTTGGGCGCGG
691	21	0	0	0	0	0	0	0	GGCTTGGGCGCGG
692	21	0	0	0	0	0	0	0	GCTTGGGCGCGG
693	21	0	0	0	0	0	0	0	CTTGGGCGCGG
694	21	0	0	0	0	0	0	0	TTGGGCGCGG
695	21	0	0	0	0	0	0	0	TGGGCGCGG
696	21	0	0	0	0	0	0	0	GGGCGCGG
697	21	0	0	0	0	0	0	0	GGCGCGG
698	21	0	0	0	0	0	0	0	GCCGCGG
699	21	0	0	0	0	0	0	0	CCGCGG
700	21	0	0	0	0	0	0	0	CGCGG
701	21	0	0	0	0	0	0	0	GCGG
702	21	0	0	0	0	0	0	0	CGG
703	21	0	0	0	0	0	0	0	GGG

111/156

704	21	0	0	0	0	0	GCGCCTCCACCTTCAAGGAGG
705	21	0	0	0	0	0	CGCTCCACCTTCAAGGAGGA
706	21	0	0	0	0	0	GCTCCACCTTCAAGGAGGAA
707	21	0	0	0	0	0	CCTCACCTTCAAGGAGGAAC
708	21	0	0	0	0	0	CTCCACCTTCAAGGAGGAACC
709	21	0	0	0	0	0	TCCACCTTCAAGGAGGAACCG
710	21	0	0	0	0	0	CCACCTTCAAGGAGGAACCGC
711	21	0	0	0	0	0	CACCTTCAAGGAGGAACCGCA
712	21	0	0	0	0	0	ACCTTCAAGGAGGAACCGCAG
713	21	0	0	0	0	0	CCTTCAAGGAGGAACCGCAGA
714	21	0	0	0	0	0	CTTCAAGGAGGAACCGCAGAC
715	21	0	0	0	0	0	TTCAAGGAGGAACCGCAGACC
716	21	0	0	0	0	0	TCAAGGAGGAACCGCAGACCG
717	21	0	0	0	0	0	CAAGGAGGAACCGCAGACCGT
718	21	0	0	0	0	0	AAGGAGGAACCGCAGACCGTG
719	21	0	0	0	0	0	AGGAGGAACCGCAGACCGTGC
720	21	0	0	0	0	0	GGAGGAACCGCAGACCGGTGCC
721	21	0	0	0	0	0	GAGGAACCGCAGACCGTGCCG
722	21	0	0	0	0	0	AGGAACCGCAGACCGTGCCGG
723	21	0	0	0	0	0	GGAACCGCAGACCGTGCCGGA
724	21	0	0	0	0	0	GAACCGCAGACCGTGCCGGAG

112/156

FIG. 24A (36)

[illegible]

113/156

FIG. 24A (37)

746	21	0	0	0	0	0	0	CGGCAGCCGGGACGCCACGC
747	21	0	0	0	0	0	0	GCGAGCCGGGACGCCACGCC
748	21	0	0	0	0	0	0	CGCAGCCGGGACGCCACGCCG
749	21	0	0	0	0	0	0	GCAGCCGGGACGCCACGCCGC
750	21	0	0	0	0	0	0	CAGCCGGGACGCCACGCCGCC
751	21	0	0	0	0	0	0	AGCCGGGACGCCACGCCGCCG
752	21	0	0	0	0	0	0	GCCGGGACGCCACGCCGCCGG
753	21	0	0	0	0	0	0	CCGGACGCCACGCCGCCCGGT
754	21	0	0	0	0	0	0	CGGACGCCACGCCGCCCGGTG
755	21	0	0	0	0	0	0	GGACGCCACGCCGCCCGGTGT
756	21	0	0	0	0	0	0	GGACGCCACGCCGCCCGGTGTC
757	21	0	0	0	0	0	0	GACGCCACGCCGCCCGGTGTCC
758	21	0	0	0	0	0	0	ACGCCACGCCGCCCGGTGTCCC
759	21	0	0	0	0	0	0	CGCCACGCCGCCCGGTGTCCCC
760	21	0	0	0	0	0	0	GCCACGCCGCCCGGTGTCCCCC
761	21	0	0	0	0	0	0	CCACGCCGCCCGGTGTCCCCCA
762	21	0	0	0	0	0	0	CACGCCGCCCGGTGTCCCCCAT
763	21	0	0	0	0	0	0	ACGCCGCCCGGTGTCCCCCATC
764	21	0	0	0	0	0	0	CGCCGCCCGGTGTCCCCCATCA
765	21	0	0	0	0	0	0	GCCGCCCGGTGTCCCCCATCAA
766	21	0	0	0	0	0	0	CCGCCCGGTGTCCCCCATCAAC

114/156

FIG. 24A (38)

767	21	0	0	0	0	0	CGCCGGTGTCCCCCATCAACA	0
768	21	0	0	0	0	0	GCCGGTGTCCCCCATCAACAT	0
769	21	0	0	0	0	0	CCGGTGTCCCCCATCAACATG	0
770	21	0	0	0	0	0	CGGTGTCCCCCATCAACATGG	0
771	21	0	0	0	0	0	GGTGTCCCCCATCAACATGGA	0
772	21	0	0	0	0	0	GTGTCCCCCATCAACATGGAA	0
773	21	0	0	0	0	0	TGTCCCCCATCAACATGGAAG	0
774	21	0	0	0	0	0	GTCCCCCATCAACATGGAAGA	0
775	21	0	0	0	0	0	TCCCCCATCAACATGGAAGAC	0
776	21	0	0	0	0	0	CCCCCATCAACATGGAAGACC	0
777	21	0	0	0	0	0	CCCATCAACATGGAAGACCA	0
778	21	0	0	0	0	0	CCCATCAACATGGAAGACCAA	0
779	21	0	0	0	0	0	CCATCAACATGGAAGACCAAG	0
780	21	0	0	0	0	0	CATCAACATGGAAGACCAAGA	0
781	21	0	0	0	0	0	ATCAACATGGAAGACCAAGAG	0
782	21	0	0	0	0	0	TCAACATGGAAGACCAAGAGC	0
783	21	0	0	0	0	0	CAACATGGAAGACCAAGAGCG	0
784	21	0	0	0	0	0	AACATGGAAGACCAAGAGCGC	0
785	21	0	0	0	0	0	ACATGGAAGACCAAGAGCGCA	0
786	21	0	0	0	0	0	CATGGAAGACCAAGAGCGCAT	0
787	21	0	0	0	0	0	ATGGAAGACCAAGAGCGCATC	0



116/156

**FIG. 24A (40)**

809	21	0	0	0	0	0	0	AAGTGGAGCGCAAGCGGCTGC
810	21	0	0	0	0	0	0	AGTGGAGCGCAAGCGGCTGCG
811	21	0	0	0	0	0	0	GTGGAGCGCAAGCGGCTGCCG
812	21	0	0	0	0	0	0	TGGAGCGCAAGCGGCTGCCGA
813	21	0	0	0	0	0	0	GGAGCGCAAGCGGCTGCCGAA
814	21	0	0	0	0	0	0	GAGCGCAAGCGGCTGCCGAAC
815	21	0	0	0	0	0	0	AGCGCAAGCGGCTGCCGAAACC
816	21	0	0	0	0	0	0	GCGCAAGCGGCTGCCGAAACCG
817	21	0	0	0	0	0	0	CGCAAGCGGCTGCCGAAACCGG
818	21	0	0	0	0	0	0	GCAAGCGGCTGCCGAAACCGGC
819	21	0	0	0	0	0	0	CAAGCGGCTGCCGAAACCGGCT
820	21	0	0	0	0	0	0	AAGCGGCTGCCGAAACCGGCTG
821	21	0	0	0	0	0	0	AGCGGCTGCCGAAACCGGCTGG
822	21	0	0	0	0	0	0	GCGGCTGCCGAAACCGGCTGGC
823	21	0	0	0	0	0	0	CGGCTGCCGAAACCGGCTGGCG
824	21	0	0	0	0	0	0	GGTGCGGAAACCGGCTGGCGCG
825	21	0	0	0	0	0	0	GCTGCGGAAACCGGCTGGCGGC
826	21	0	0	0	0	0	0	CTGCGGAACCGGCTGGCGGCC
827	21	0	0	0	0	0	0	TGCGGAACCGGCTGGCGGCCA
828	21	0	0	0	0	0	0	GCGGAACCGGCTGGCGGCCAC
829	21	0	0	0	0	0	0	CGGAACCGGCTGGCGGCCACC





118/156

FIG. 24A (42)

851	21	0	0	0	0	0	AGTCCGGAAGCGGAAGCTGG
852	21	0	0	0	0	0	GTGCCGGAAGCGGAAGCTGGA
853	21	0	0	0	0	0	TGCCGGAAGCGGAAGCTGGAG
854	21	0	0	0	0	0	GCCGGAAGCGGAAGCTGGAGC
855	21	0	0	0	0	0	CCGGAAGCGGAAGCTGGAGCG
856	21	0	0	0	0	0	CGAAGCGGAAGCTGGAGCGC
857	21	0	0	0	0	0	GGAAGCGGAAGCTGGAGCGCA
858	21	0	0	0	0	0	GAAGCGGAAGCTGGAGCGCAT
859	21	0	0	0	0	0	AAGCGGAAGCTGGAGCGCATC
860	21	0	0	0	0	0	AGCGGAAGCTGGAGCGCATCG
861	21	0	0	0	0	0	GCGGAAGCTGGAGCGCATCGC
862	21	0	0	0	0	0	CGGAAGCTGGAGCGCATCGCG
863	21	0	0	0	0	0	GGAAGCTGGAGCGCATCGCGC
864	21	0	0	0	0	0	GAAGCTGGAGCGCATCGCGCG
865	21	0	0	0	0	0	AAGCTGGAGCGCATCGCGCGC
866	21	0	0	0	0	0	AGCTGGAGCGCATCGCGCGCC
867	21	0	0	0	0	0	GCTGGAGCGCATCGCGCGCCT
868	21	0	0	0	0	0	CTGGAGCGCATCGCGCGCCTG
869	21	0	0	0	0	0	TGGAGCGCATCGCGCGCCTGG
870	21	0	0	0	0	0	GGAGCGCATCGCGCGCCTGGA
871	21	0	0	0	0	0	GAGCGCATCGCGCGCCTGGAG

119/156

**FIG. 24A (43)**

[illegible]

120/156

**FIG. 24A (44)**

[illegible]

121/156

FIG. 24A (45)

914	21	0	0	0	0	0	CCGAGAACGGCGGGCTGTCTCGA
915	21	0	0	0	0	0	CGAGAACGGCGGGCTGTCTCGAG
916	21	0	0	0	0	0	GAGAACGGCGGGCTGTCTCGAGT
917	21	0	0	0	0	0	AGAAACGGCGGGCTGTCTCGAGTA
918	21	0	0	0	0	0	GAACGCCGGGGCTGTCTCGAGTAC
919	21	0	0	0	0	0	AACGCCGGGGCTGTCTCGAGTACC
920	21	0	0	0	0	0	ACGCCGGGGCTGTCTCGAGTACCG
921	21	0	0	0	0	0	CGCCGGGGCTGTCTCGAGTACCGC
922	21	0	0	0	0	0	GCGGGCTGTCTCGAGTACCGCCC
923	21	0	0	0	0	0	CGGGCTGTCTCGAGTACCGCCCG
924	21	0	0	0	0	0	GGGCTGTCTCGAGTACCGCCCGG
925	21	0	0	0	0	0	GGCTGTCTCGAGTACCGCCCGGC
926	21	0	0	0	0	0	GGCTGTCTCGAGTACCGCCCGGCC
927	21	0	0	0	0	0	GCTGTCGAGTACCGCCCGGCCCT
928	21	0	0	0	0	0	CTGTCGAGTACCGCCCGGCCCTC
929	21	0	0	0	0	0	TGTCGAGTACCGCCCGGCCCTCC
930	21	0	0	0	0	0	GTCGAGTACCGCCCGGCCCTCCT
931	21	0	0	0	0	0	TCCGAGTACCGCCCGGCCCTCCTC
932	21	0	0	0	0	0	CGAGTACCGCCCGGCCCTCCTCC
933	21	0	0	0	0	0	GAGTACCGCCCGGCCCTCCTCCG
934	21	0	0	0	0	0	AGTACCGCCCGGCCCTCCTCCGG

122/156

FIG. 24A (46)

935	21	0	0	0	0	0	GTACCGCCGGCCTCCTCCGCGG
936	21	0	0	0	0	0	TACCGCCGGCCTCCTCCGCGGA
937	21	0	0	0	0	0	ACGCCGGCCTCCTCCGCGGAG
938	21	0	0	0	0	0	CCGCCGCCTCCTCCGCGGAGC
939	21	0	0	0	0	0	CGCCGGCCTCCTCCGCGGAGCA
940	21	0	0	0	0	0	GCCGGCCTCCTCCGCGGAGCAG
941	21	0	0	0	0	0	CCGGCCTCCTCCGCGGAGCAGG
942	21	0	0	0	0	0	CGCCTCCTCCGGGAGCAGGT
943	21	0	0	0	0	0	GGCCTCCTCCGGGAGCAGGTG
944	21	0	0	0	0	0	GCCTCCTCCGGGAGCAGGTGG
945	21	0	0	0	0	0	CCTCCTCCGGGAGCAGGTGGC
946	21	0	0	0	0	0	CTCCTCCGGGAGCAGGTGGCC
947	21	0	0	0	0	0	TCCTCCGGGAGCAGGTGGCCC
948	21	0	0	0	0	0	CCTCCGGGAGCAGGTGGCCCCA
949	21	0	0	0	0	0	CTCCGGGAGCAGGTGGCCCCAG
950	21	0	0	0	0	0	TCCGGGAGCAGGTGGCCCCAGC
951	21	0	0	0	0	0	CCGGGAGCAGGTGGCCCCAGCT
952	21	0	0	0	0	0	CGGAGCAGGTGGCCCCAGCTC
953	21	0	0	0	0	0	GGGAGCAGGTGGCCCCAGCTCA
954	21	0	0	0	0	0	GGAGCAGGTGGCCCCAGCTCAA
955	21	0	0	0	0	0	GAGCAGGTGGCCCCAGCTCAAA

123/156

FIG. 24A (47)

956	21	0	0	0	0	0	AGCAGGTGGCCCAGCTCAAAAC
957	21	0	0	0	0	0	GCAGGTGCCCCAGCTCAAACA
958	21	0	0	0	0	0	CAGGTGGCCCAGCTCAAACAG
959	21	0	0	0	0	0	AGGTGGCCCAGCTCAAAACAGA
960	21	0	0	0	0	0	GGTGGCCCAGCTCAAACAGAAA
961	21	0	0	0	0	0	GTGGCCAGCTCAAACAGGAAG
962	21	0	0	0	0	0	TGGCCCAGCTCAAACAGGAAGG
963	21	0	0	0	0	0	GGCCAGCTCAAACAGGAAGGT
964	21	0	0	0	0	0	GCCAGCTCAAACAGGAAGGTC
965	21	0	0	0	0	0	CCAGCTCAAACAGGAAGGTCA
966	21	0	0	0	0	0	CCAGCTCAAACAGGAAGGTCAT
967	21	0	0	0	0	0	CAGCTCAAACAGGAAGGTCATG
968	21	0	0	0	0	0	AGCTCAAACAGGAAGGTCATGA
969	21	0	0	0	0	0	GCTCAAACAGGAAGGTCATGAC
970	21	0	0	0	0	0	CTCAAACAGGAAGGTCATGACC
971	21	0	0	0	0	0	TCAAACAGGAAGGTCATGACCC
972	21	0	0	0	0	0	CAAACAGGAAGGTCATGACCCA
973	21	0	0	0	0	0	A AACAGGAAGGTCATGACCCAC
974	21	0	0	0	0	0	AACAGGAAGGTCATGACCCACG
975	21	0	0	0	0	0	ACAGGAAGGTCATGACCCACGT
976	21	0	0	0	0	0	CAGGAAGGTCATGACCCACGTC

124/156

[illegible]





FIG. 24A (50)

1016 21 0 0 0 0 0 TGCTTGGGGTCAAGGGACACG  
1017 21 0 0 0 0 0 GCTTGGGGTCAAGGGACACGC  
1018 21 0 0 0 0 0 CTTGGGGTCAAGGGACACGCC  
1019 21 0 0 0 0 0 TTGGGGTCAAGGGACACGCCT  
1020 21 0 0 0 0 0 TGGGGTCAAGGGACACGCCTT  
1021 21 0 0 0 0 0 GGGTCAAGGGACACGCCCTTC  
1022 21 0 0 0 0 0 GGGTCAAGGGACACGCCCTTCT  
1023 21 0 0 0 0 0 GGTCAGGGACACGCCCTTCTG  
1024 21 0 0 0 0 0 GTCAAGGGACACGCCCTTCTGA

126/156

FIG. 24B (1)

(Partial File -- 10 pages of 190 pages)

OligoProbe DesignStation

Probes: C:\HITACHI\HUMBUNX.CDS  
Preparation: C:\HITACHI\JUNMIX.PRP

127/156

	Locus pos	Tm	Locus pos	Tm	Locus pos	Tm
atgtgcactaaaatgggaacagcccttctac	1	30	1	1	2	2
humbjunc 1	1	60.76	2	2	3	4
musbjunc 1	1	50.03				
muscjunc 1	1	30.07				
musdjunc 721	721	27.84				

FIG. 24B (2)

tgtgcactaaaatgggaacagcccttctac  
 2 29 1 1 2  
 humbjunx 65533 60.68  
 musbjunx 65533 49.58  
 muscjunx 1 29.97  
 musdjunx 721 27.66

gtgcactaaaatgggaacagcccttctac  
 3 28 1 1 2  
 humbjunx 65533 60.60  
 musbjunx 65533 49.10  
 muscjunx 1 29.86  
 musdjunx 721 27.47

128/156

FIG. 24B (3)

tgactataaatggaacagcccttctacc  
 4 28 1 1 1 1 2 2 2 2 3 4  
 humbjunx 65533 60.60  
 musbjunx 65533 46.57  
 muscjunx 1 29.86  
 musdjunx 729 27.47

129/156

gcactataaatggaacagcccttctacc  
 5 27 1 1 1 1 2 2 2 2 3 4  
 humbjunx 5 60.51  
 musbjunx 5 45.96  
 muscjunx 1 29.75  
 musdjunx 729 27.26

130/156

ctaaaatggaacagcccttctaccag	1	2	2	3	4
8 27	1	1			
humbjunc	1	60.51			
musbjunc	5	45.96			
muscjunc	1	33.33			
musdjunc	729	27.26			

FIG. 24B (5)

131/156

<b>aaaaatggaacagcccttctaccacgc</b>							
10	27	1	1	2	2	3	4
humbjux	5	60.51					
musbjux	5	49.70					
muscjux	9	34.44					
musdjux	729	27.26					

[illegible]

FIG. 24B (6)

132/156

aatggaacagcccttctaccacgac									
12	25	1	1	2	2	2	3	3	4
humbj	j	unx	5						
60.32									
musbj	j	unx	5						
48.64									
muscj	j	unx	9						
34.56									
musdj	j	unx	737						
26.80									
atggaacagcccttctaccacgac									
13	24	1	1	2	2	2	4	5	6
humbj	j	unx	13						
60.20									
musbj	j	unx	13						
48.04									
muscj	j	unx	9						
34.62									
musdj	j	unx	1						
32.46									
humdj	j	unx	65533						
30.25									
musdj	j	unx	737						
26.55									

FIG. 24B (7)

133/156

tggaacagcccttctaccacgac	1	2	2	2	3	5	6
14 23 1 1	1						
humbjunc 9	60.08						
musbjunc 9	47.39						
muscjunc 9	33.39						
musdjunc 1	31.14						
humdjunc 65533	28.83						
musdjunc 737	26.27						
ggaacagcccttctaccacgac	2	2	2	2	3	5	6
15 23 1 1	1						
humbjunc 9	61.86						
musbjunc 9	49.17						
muscjunc 9	32.09						
musdjunc 1	29.83						
humdjunc 65533	28.53						
musdjunc 737	26.27						



134/156

FIG. 24B (8)

gaacagcccttctaccacgacga	1	2	2	2	4	8
16 23 1 1	musdjunx	737	26.27			
humbjunx 9						
musbjunx 9						
muscjunx 9						
humdjunx 65533						
musdjunx 1						
humbjunx 281						
musbjunx 281						
17 23 1 1	musdjunx	737	26.27	6		8
humbjunx 17						
musbjunx 17						
muscjunx 17						
humdjunx 5						
humbjunx 281						
musbjunx 281						
musdjunx 1						

FIG. 24B (9)

[illegible]

136/156

FIG. 24B (10)

agcccttctaccacgacgactca	1	2	2	2	6	7
20 23	1					
humbjunc	13					
musbjunc	17					
muscjunc	17					
humdjunc	5					
humbjunc	281					
musbjunc	281					
musdjunc	1					
gcccttctaccacgacgactcat	1	2	2	2	6	7
21 23	1					
humbjunc	21					
musbjunc	17					
muscjunc	17					
humdjunc	5					
humbjunc	281					
musbjunc	281					
musdjunc	9					

137/156

FIG. 24B (11)

cccttctaccacgacgactcatcac	1	1	2	2	6	7
22 24 1 1						
humbjunc 17						
musbjunc 17						
humbjunc 281						
muscjunc 17						
humdjunc 5						
musbjunc 281						
musdjunc 5						
cccttctaccacgacgactcatcac	1	1	2	2	5	6
23 25 1 1						
humbjunc 17						
musbjunc 17						
humbjunc 289						
muscjunc 17						
musbjunc 289						
humdjunc 5						

**FIG. 24B (12)**

cttctaccagcagctcatacacag	1	2	2	3	4
24 26 1 1	1	2	2	3	4
humbjux 17	60.42				
musbjux 17	44.00				
humbjux 289	35.65				
musbjux 289	29.77				

[illegible]

tctaccacgacgactcatacacagc					
26	25	1	1	1	
humbjunc	21				
musbjunc	25				
humbjunc	289				
musbjunc	289				

FIG. 24B (13)

ctaccagcagcactcatacacagc	1	2	2	2	3	4	4
27 24 1 1 1							
humbjunc 21 60.20							
musbjunc 25 45.37							
humbjunc 289 35.87							
musbjunc 289 29.50							

taccagcagcactcatacacagctac	1	1	2	2	3	4	4
28 26 1 1 1							
humbjunc 21 60.42							
musbjunc 25 42.26							
humbjunc 289 35.65							
musbjunc 289 29.77							

accagcagcactcatacacagctac	1	1	2	2	3	4	4
29 25 1 1 1							
humbjunc 29 60.32							
musbjunc 25 42.64							
humbjunc 289 35.76							
musbjunc 289 29.64							

139/156

FIG. 24B (14)

ccagcagactcatacacagctac									
30 24	1	1	1	1	1	1	1	1	1
humbjunc	25								
musbjunc	25								
humbjunc	289								
musbjunc	289								

cacgacgactcatacacagctacg									
31 24	1	1	1	1	1	1	1	1	1
humbjunc	25								
musbjunc	25								
humbjunc	297								
musbjunc	297								
humdjunc	573								

acgacgactcatacacagctacgg									
32 24	1	1	1	1	1	1	1	1	1
humbjunc	25								
musbjunc	25								
humbjunc	293								
humdjunc	573								
musbjunc	293								

140/156

141/156

**FIG. 24B (15)**

[illegible]

gacgactacatacacagctacggg									
34	23	1	1	1	1	1	2	2	3
humbj	j	unx	29	60.08					
musbj	j	unx	29	41.82					
humdj	j	unx	581	26.27					

acgactcatacacagctacgggatac									
35	26	1	1	1	2	2	2	2	3
humbj	j	29	60.42						
musbj	j	29	44.26						
humdj	j	581	27.04						



142/156

actc	taca	cagc	ctac	ggga	tacg	gata	cg
38	25	1	1	1	1	1	1
humbj	jux	33		60.32			
musbj	jux	37		43.52			
humdj	jux	581		26.80			

FIG. 24B (17)

143/156

ctc	1	1	2	2	2	3
tac						
agg						
gata						
cgg						
39	1	1	1	1	1	1
24	1	1	1	1	1	1
humbj	33	60.20				
musbj	37	42.70				
humdj	581	26.55				

tc	1	1	1	2	2	3
ata						
cag						
ctac						
gggata						
cgc						
40	1	1	1	1	1	1
24	1	1	1	1	1	1
humbj	33	60.20				
musbj	37	39.75				
humdj	581	26.55				

cat	1	1	1	2	2	3
ata						
cag						
ctac						
gggata						
cgc						
41	1	1	1	1	1	1
23	1	1	1	1	1	1
humbj	41	60.08				
musbj	37	38.91				
humdj	581	26.27				

FIG. 24B (18)

atacacagctacgggatacggcc	1	1	1	2	2	2	3
42 23 1 1							
humbjnx 37							
musbjnx 37							
humdjnx 589							
							26.27

tacacagctacgggatacggccg	1	1	2	2	2	2	3
43 23 1 1							
humbjnx 37							
musbjnx 37							
humdjnx 589							
							26.27

acacagctacgggatacggccg	1	1	2	2	2	2	3
44 22 1 1							
humbjnx 37							
musbjnx 37							
humdjnx 589							
							25.96

144/156

FIG. 24B (19)

145/156

cacagctacgggatacggccg	1	1	2	2	2	3
45 21 1	1					
humbjux 45	61.76					
musbjux 45	40.00					
humdjux 589	25.62					

acagctacgggatacggccg	1	1	2	2	2	3
46 21 1	1					
humbjux 41	61.76					
musbjux 45	43.38					
humdjux 589	25.62					

FIG. 24B (20)

cagctacgggatacggccgg  
 47 20 1 1 1 1 1 1 1  
 humbjunx 41 61.70  
 musbjunx 45 43.90  
 humdjunx 589 25.25

agctacgggatacggccgg  
 48 20 1 1 1 1 1 1 1  
 humbjunx 41 61.70  
 musbjunx 45 40.35  
 muscjunx 561 31.40

146/156

147/156

FIG. 25

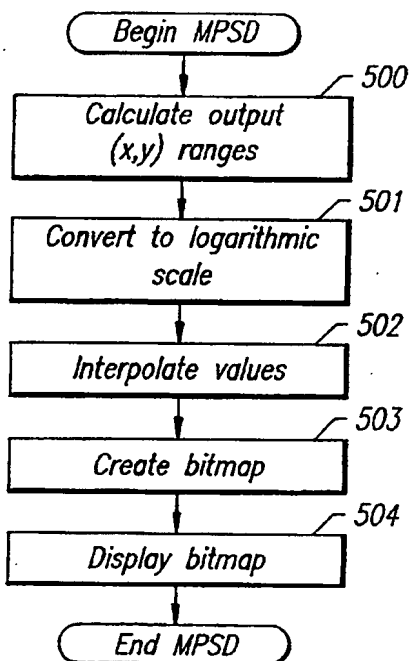
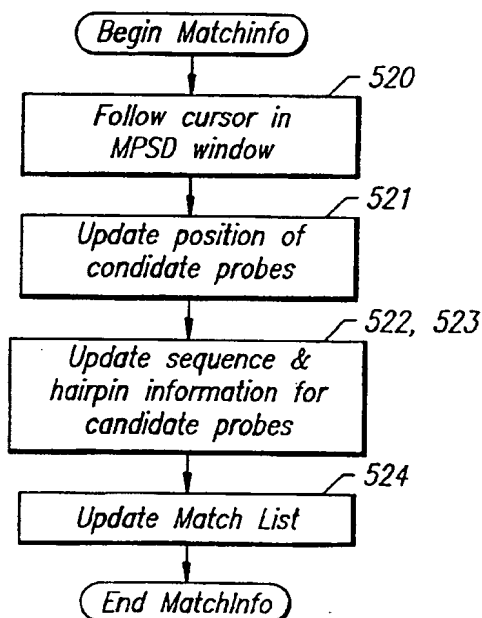


FIG. 26



148/156

## FIG. 27

LOCUS HUMBUNX 1044 bp DNA 19-DEC-1991  
 BASE COUNT 195 A 368 C 340 G 141 T  
 ORIGIN

```

1 ATGTGCACTA AAATGGAACA GCCCTTCTAC CACGACGACT CATAACAGC TACGGGATAC
61 GGCCGGGCCC CTGGTGGCCT CTCTCTACAC GACTACAAAC TCCTGAAACC GAGCCTGGCG
121 GTCAACCTGG CCGACCCCTA CCGAGTCTC AAAGCGCCTG GGGTCGCGG ACCCGGCCCA
181 GAGGGCGGCG GTGGCGGCAG CTACTTTTCT GGTACGGCT CGGACACCG CGGTCCTCTC
241 AAGCTCGCCT CTTCGGAGCT GGAACGCCGT ATGTCCCCA ACAGCAACGG CGTGATCACG
301 ACGACGCCTA CACCCCGGG ACAGTACTTT TACCCCGCG GGGTGGCAG CGGTGGAGGT
361 GCAGGGGGCG CAGGGGGCGG CGTCACCGAG GAGCAGGAGG GCTTCGCCGA CGGCTTTGTC
421 AAAGCCCTGG ACGATCTGCA CAAGATGAAC CACGTGACAC CCCCACACGT GTCCCTGGGC
481 GCTACCGGGG GCGCCCGCG CTACTCCCCA GCCTCTGCGT CCTCGGAGG GCCACCTCCC
541 GTTTACACCA ACCTCAGCAG GTAGCCCGACG ACCACCATCA GGTACCTCCC ACACGCGCCG
601 GCCGTCGGGA CCGGGAGCTC GTACCCGACG GCGCAGCTG GGCTTGGCC GCGGCGCCTC CACCTTCAAG
661 CCTTTCGCCG GTGGCCACCC GGCGAGTGC GGAGGCGCG CCGGCGCGG CCACGCGGCC GGTGTCCCCC
721 GAGGAACCGC AGACCGTGCC GGAGGCGCG CCGCATCAA GTGGAGCGCA AGCGGTGCG GAACCGGCTG
781 ATCAACATGG AAGACCAAGA GCGCAAGCTG GAGCGCATCG CGGCGCTGGA GGACAAGGTG
841 GCGGCCACCA AGTCCCGGAA GCGGAAGCTG TCGAGTACCG CCGGCTCCT CCGGAGCAG
901 AAGACGCTCA AGGCCGAGAA GCGGGGCTG CACGTCAGCA ACGGCTGTCA GCTGCTGCTT
961 GTGGCCCCAGC TCAACACAGAA GGTATGACC
1021 GGGGTCAAGG GACACGCCTT CTGA
  
```

//

FIG. 28 (1)

LOCUS HUMBUNX 1044 bp DNA 19-DEC-1991  
 BASE COUNT 195 A 340 G 141 T  
 ORIGIN

1 ATGTGCACTA AAATGGAACA GCCCTTCTAC CACGACGACT CACACACAGC TACGGGATAC  
 61 GGCCGGGCCC CTGGTGGCCT CTCTCTACAC GACTACAAAC TCCTGAAACC GAGCCTGGCG  
 121 GTCAACCTGG CCGACCCCTA CCGAGTCTC AAAGCGCCTG GGGCTCGCGG ACCCGGCCCA  
 181 GAGGCGGCG GTGGCGGCG CTACTTTCT GGTACGGGT ACAGCAACCG CGGTGAGGT  
 241 AAGCTCGCT CTTCGGAGT GGAACGCTG ATTGTCCCA TACCCCGCG GCTTCGCCA CGGCTTTGTC  
 301 ACGACGCCTA CACCCCGGG ACAGTACTT GAGCAGGAG CCGGCGGAG GTCCCTGGG  
 361 GCAGGGGCG CAGGGGCGG CGTCACCGAG CACGTGACAC CCGGCCCGG GCCACCTCCC  
 421 AAGCCCTGG ACGATCTGCA CAAGATGAAC GGCGTCTACG CCTCGGAGG ACACGCGCG  
 481 GCTACCGGG GGCCCGCGG TGGCCCGGG CTACTCCCA GCCTCTGCGT GCGGCGGCT  
 541 GTTACACCA ACCTCAGCAG GTACCCGACG ACCACCATCA GGCTTGGCC AGCGCGCTC  
 601 CCGTCGGA CCGGAGCTC GTGCGAGCTG GCGCGAGCTG AGCGGGAGC GTGAGCGCA  
 661 CCCTTCGCG AGACCGTGCC GGAGCGCGC GCGCATCAA GTGAGCGCA AGCGGCTGCG  
 721 GAGGAACCG AGACCGTGCC GGAGCGCGC GCGCATCAA GTGAGCGCA AGCGGCTGCG  
 781 ATCAACATGG AAGACCAAGA GCGGAGCTG GAGCGCATCG CGGCGCTGGA GAACCGGCTG  
 841 GCGGCCACCA AGTGCCGGAA GCGGAGCTG TCGAGTACCG CCGGCTCCT CCGGAGCAG  
 901 AAGACGCTCA AGGCGAGAA CGCGGGGCTG CACGTACGCA ACGGCTGTCA GCTGCTGCTT  
 961 GTGGCCCGAG TCAAACAGAA GGTATGACC CACGTACGCA ACGGCTGTCA GCTGCTGCTT  
 1021 GGGGTCAAGG GACACGCCCTT CTGA

//



FIG. 28 (2)

LOCUS HUMCJUNX 996 bp DNA 19-DEC-1991  
 BASE COUNT 226 A 342 C 299 G 129 T  
 ORIGIN

1 ATGACTGCAA AGATGGAAC GACCTTCTAT GACGATGCCC TCAACGCCTC GTTCCTCCCG  
 61 TCCGAGAGCG GACCTTATGG CTACAGTAAC CCCAAGATCC TGAAACAGAG CATGACCCCTG  
 121 AACCTGGCCG ACCCAGTGGG GAGCCTGAAG CCGCACCTCC CGCCCAAGAA CTCGGACCTC  
 181 CTCACCTCGC CCGACGTGGG GCTGCTCAAG CTGGCGTCGC CCGAGCTGGA GCGCCTGATA  
 241 ATCCAGTCCA GCAACGGGCA CATCACCACC ACGCCGACCC CCACCCAGTT CCTGTGCCCC  
 301 AAGAACGTGA CAGATGAGCA GGAGGGGTTT GCCAGGGCT TCGTGCGCGC CCTGGCCGAA  
 361 CTGCACAGCC AGAACACGCT GCCAGCGTC ACGTCGGCGG CGCAGCCGGT CAAAGGGGCA  
 421 GGCAATGGTG CTCCCGCGGT AGCCTCGGTG GCAGGGGCA GCGCAGCGG CGGCTTCAGC  
 481 GCCAGCCTGC ACAGCGAGCC GCCGGTCTAC GCAAACCTCA GCCTGGCCTT TCCCGCGCAA  
 541 CTGAGCAGCG GCGCGGGGC GCCCTCCTAC GGCGCGCGCG AGATGCCCGT GCAGACCCG  
 601 CCCAGCAGC AGCAGCAGCC GCCGACCCAC CTGCCCCCAGC ACATGCCCGG AGAGACACCG  
 661 CGGCTGCAGG CCTGAAGGA GGAGCCTCAG ACAGTGCCCG GAGCGGATCA AGGCGGAGAG  
 721 CCCCTGTCCC CCATCGACAT GGAGTCCCAG GAGCGGATCA AAAAGGAAGC TGGAGAGAAT  
 781 AGGAACCGCA TCGCTGCCCTC CAAGTGCCGA AACTCGGAGC TGGCGTCCAC GCGCCGGCTG  
 841 GAGGAAAAAG TGAAAAACCTT GAAAGCTCAG AACTCGGAGC TGGCGTCCAC GGCCAAACATG  
 901 CTCAGGGAAC AGGTGGCACA GCTTAAACAG AAAGTCATGA ACCACGTTAA CAGTGGGTGC  
 961 CAACTCATGC TAACGCAGCA GTTGCAACA TTTTGA

150/156

151/156

## FIG. 28 (3)

LOCUS	HUMDJUNX	1044 bp ss-mRNA	PRI	24-MAY-1991
DEFINITION	Human junD mRNA			
ACCESSION	X56681			
KEYWORDS	jun-D gene; oncogene.			
SOURCE	Homo sapiens RNA.			
ORGANISM	Homo sapiens			
	Eukaryota; Animalia; Metazoa; Chordata; Vertebrata; Mammalia;			
	Theria; Eutheria; Primates; Haplorhini; Catarrhini; Homnidae.			
REFERENCE	1 (bases 1 to 1891)			
AUTHORS	Shaul, Y.			
JOURNAL	Unpublished (1990)			
STANDARD	full automatic			
REFERENCE	2 (sites)			

152/156

## FIG. 28 (4)

AUTHORS Berger, I. and Shaul, Y.  
 TITLE Structure and function of human jun-D  
 JOURNAL Unpublished (1990)  
 STANDARD full staff review  
 COMMENT From EMBL 26 entry HSJUNDR; dated 18-MAR-1991.  
 FEATURES  
     mRNA  
         Location/Qualifiers  
         1..1891  
         /gene="junD"  
         /evidence=EXPERIMENTAL  
 CDS  
     175..1218  
     /product="junD protein"  
     /gene="junD"  
     /codon\_start=1  
 polyA\_site  
     1891..1891

FIG. 28 (5)

BASE COUNT            162 A            405 C            360 G            117 T  
 ORIGIN  
 1 ATGGAACAC CCTTCTACCG CGATGAGGCG CTGAGCGGCC TGGCGGGCGG CGCCAGTGGC  
 61 AGCGGGGCA CGTTCGCGTC CCCGGCCGC TTGTTCCCG GGGCGCCCC GACGGCCGG  
 121 GCCGGCAGCA TGATGAAGAA GGACGCGCTG ACGCTGAGCC TGAGTGAGCA GGTGGCGGCA  
 181 GCGCTCAAGC CTGCGCCCCG GCGCGCCTCC TACCCCTG CCGCGACGG CGCCCCCAGC  
 241 GCGCACCCC CCGACGGCCT GCTCGCCTCT CCCGACCTGG GGCTGCTGAA GCTGGCCTCC  
 301 CCCGAGCTCG AGCGCTCAT CATCCAGTCC AACGGCTGG TCACCAACCAC GCCGACGAGC  
 361 TCACAGTTCC TCTACCCCAA GGTGGCGGCC AGCGAGGAGC AGGAGTTGCG CGAGGGCTTC  
 421 GTCAAGGCCC TGGAGGATT ACACAAGCAG AACAGCTCG GCGGGGCGG GCGCGCTGCC  
 481 GCCGCCGCCG CCGCCGCCG GGGCCCTCG GGCAAGGCCA CCGGCTCCG GCGCGCTGCC  
 541 GAGCTGGCCC CCGCGCGGC CGGCCCGAA GCGCCTGTCT ACGCGAACCT GAGCAGCTAC  
 601 GCGGGCGGCG CCGGGGCGC GGGGGCGCC GCGACGGTCG CCTTCGCTGC CGAACCTGTG  
 661 CCTTCCCGC CGCCGCCAC AGACGGTGCC CCCAGGCGG TTGGGGCCG CGGCTTGGC TGGCTCAAG  
 721 GACGAGCCAC AGACGGTGCC CGACGTGCG AGCTTCGGG AGAGCCCGC GTTGTGCCC  
 781 ATCGACATGG ACAGCAGGA GCGATCAAG GCGGAGCGCA AGCGGCTGCG CAACCGCATC  
 841 GCCGCCTCCA AGTGCCGCAA GCGCAAGCTG GAGCGCATCT CCGCCCTGGA AGAGAAAGTG  
 901 AGACCCCTCA AGAGTCAGAA CACGGAGCTG GGTCCACGG CGAGCCTGCT GCGGAGCAG  
 961 GTGGCGCAGC TCAAGCAGAA AGTCCTCAGC CACGTCAACA GCGGCTGCCA GCTGCTGCCC  
 1021 CAGCACCCAG TCCCGGCGTA CTGA  
 //

153/156

FIG. 28 (6)

LOCUS	MUSBJUNX	1035 bp	DNA	
BASE COUNT	210 A	333 C	333 G	159 T
ORIGIN				19-DEC-1991

1	ATGTGCACGA	AAATGGAACA	GCCTTTCTAT	CACGACGACT	CTTACGCAGC	GGCGGGATAC
61	GGTCGGAGCC	CTGGCAGCCT	GTCTCTACAC	GACTACAAAC	TCCTGAAACC	CACCTTGGCG
121	CTCAACCTGG	CGGATCCCTA	TCGGGGTCTC	AAGGGTCCTG	GGCGCGGGG	TCCAGGCCCG
181	GAGGGCAGTG	GGCAGGCAG	CTACTTTTCG	GGTCAGGGAT	CAGACACAGG	CGCATCTCTG
241	AAGCTAGCCT	CCACGGAAC	GGAGCGCTTG	ATCGTCCCCA	ACAGCAACGG	CGTGATCACG
301	ACGACGCCCC	CGCCTCCGGG	ACAGTACTTT	TACCCCCGTG	GGGTGGGCAG	CGGTGGAGGT
361	ACAGGGGGCG	GCGTCACCGA	GGAGCAGGAG	GGCTTTGCGG	ACGGTTTGT	CAAAGCCCCTG
421	GACGACCTGC	ACAAGATGAA	CCACGTGACG	CCCCCCAACG	TGTCCCTGGG	CGCCAGCCGGG
481	GGTCCCCCAGG	CCGGCCCCAGG	GGCGTCTAT	GCTGGTCCGG	AGCCGCCCTCC	CGTCTACACC
541	AACCTCAGCA	GTTACTCTCC	AGCCTCTGCA	CCCTCTGGAG	GCTCCGGGAC	CGCCGTCCGG
601	ACTGGGAGCT	CATACCCGAC	GGCCACCATC	AGCTACCTCC	CACATGCACC	ACCCTTTGGC
661	GGCGGCCACC	CGGCACAGCT	GGGTTGAGT	CGTGGCGCTT	CCGCCTTAA	AGAGGAACCG
721	CAGACCGTAC	CGGAGGCACG	CAGCCGCGAC	GCCACGCCGC	CTGTGTCCCC	CATCAACATG
781	GAAGACCAGG	AGCGCATCAA	AGTGGAGCGA	AAGCGGCTGC	GGAACAGGCT	GGCGGCCACC
841	AAGTGCCGGA	AGCGGAAGCT	GGAGCGCATC	GCGCGCCTGG	AGGACAAAGT	GAAGACACTC
901	AAGGCTGAGA	ACGCGGGGCT	GTCGAGTGCT	GCCGGTCTCC	TAAGGGAGCA	AGTGGCGCAG
961	CTCAAGCAGA	AGGTCATGAC	CCATGTCAGC	AACGGCTGCC	AGTTGCTGCT	AGGGGTCAAG
1021	GGACACGCCT	TCTGA				

//

FIG. 28 (7)

155/156

LOCUS	MUSCJUNX	1005 bp	DNA	19-DEC-1991		
BASE COUNT	223 A	334 C	300 G			
ORIGIN			148 T			
1	ATGACTGCAA	AGATGGAAC	GACCTTCTAC	TCAACGCCTC	GTTCTCCAG	
61	TCCGAGAGCG	GTGCCTACGG	CTACAGTAAC	CCTAAGATCC	TAAACACAGAG	CATGACCTTG
121	AACCTGGCCG	ACCCGGTGGG	CAGTCTGAAG	CCGACCTCC	GCGCAAGAA	CTCGGACCTT
181	CTCACGTCGC	CCGACGTCGG	GCTGCTCAAG	CTGCGTCGC	CGGAGCTGGA	GCGCTGATC
241	ATCCAGTCCA	GCAATGGGA	CATCACCACT	ACACCGACCC	CCACCCAGTT	CTTGTGCCCC
301	AAGAACGTGA	CCGACGAGCA	GGAGGCTTC	GCCGAGGGCT	TCGTGCGCGC	CCTGGCTGAA
361	CTGCATAGCC	AGAACACGCT	TCCAGTGTC	ACCTCCGCGG	CACAGCCGGT	CAGCGGGCGG
421	GGCATGGTGG	CTCCCGCGGT	GGCCTCAGTA	GCAGGCGCTG	GCGGCGGTGG	TGGCTACAGC
481	GCCAGCCTGC	ACAGTGAGCC	TCCGGTCTAC	GCCAACTCA	GCAACTTCAA	CCCGGTGCG
541	CTGAGCAGCG	GCGTGCGGC	GCCCTCCTAT	GCGCGGCGG	GGCTGGCCTT	TCCCTCGCAG
601	CCGCAGCAGC	AGCAGCAGCC	GCCTCAGCCG	CCGCACCACT	TGCCCCAACA	GATCCCCGTG
661	CAGCACCCGC	GGCTGCAAGC	CCTGAAGGAA	GAGCCGCAGA	CCGTGCCGGA	GATGCCGGA
721	GAGACGCCGC	CCCTGTCCCC	TATCGACATG	GAGTCTCAGG	AGCGGATCAA	GGCAGAGAGG
781	AAGCGCATGA	GGAAACCGCAT	TGCCGCCTCC	AAGTGCCGGA	AAAGGAAAGCT	GGAGCGGATC
841	GCTCGGCTAG	AGGAAAAAGT	GAAAACCTTG	AAAGCGCAA	ACTCCGAGCT	GGCATCCACG
901	GCCAACATGC	TCAGGGAACA	GGTGGCACAG	CTTAAGCAGA	AAGTCATGAA	CCACGTTAAC
961	AGTGGGTGCC	AACTCATGCT	AACGCAGCAG	TTGCAAAACGT	TTTGA	

//

FIG. 28 (8)

LOCUS	MUSDJUNX	1026 bp	DNA	19-DEC-1991
BASE COUNT	172 A	382 C	343 G	
ORIGIN				
1	ATGGAACGC	CCTTCTATGG	CGAGGAGGCG	CTGAGCGGCC
61	GTGCTGGTG	CTACTGGGGC	CCCCGGCGGT	GGTGGCTTCG
121	CCCGGGGCG	CCCCGACGAG	CAGCATGCTG	AAGAAAGACG
181	GAGCAGGAG	CGGCGGGATT	GAAACCAAGG	TCGGCCACTG
241	GACGGCGCC	CCGACGGGCT	GCTGGCTTCG	CCGGATCTTG
301	CCGGAGCTG	AGAGGCTGAT	CATCCAGTCC	AACGGGCTGG
361	ACGCAGTTC	TCTACCCGAA	GGTGGCAGCC	AGCGAGGAGC
421	GTCAAGGCG	TGGAGGACCT	GCACAAGCAA	AGCCAGCTGG
481	TCAGGGGCTC	CCGCGCCTCC	CGCGCCCGCC	GACCTGGCCG
541	ACCCCGGTCT	ACGCCAACCT	GAGCAGTTTC	GCGGTGGCG
601	GCCACCGTGG	CTTTCGCCGC	GGAGCCAGTG	CCCTTCCCGC
661	CCGCCGCCAC	CTCCGCATCC	ACCGCGCCTG	GCCGCGCTCA
721	CCGACCGTGC	CGAGCTTCGG	CGACAGCCCT	CCGCTGTCGC
781	GAACGCATCA	AGGCGGAGCG	CAAGAGGCTG	CGCAACCGCA
841	AAGCGCAAGC	TGGAGCGTAT	CTCGCGCCTG	GAGGAGAAAG
901	AACACCGAGC	TGGCGTCCAC	CGCCAGCCTG	CTGCGCGAGC
961	AAAGTCCTCA	GCCACGTCAA	CAGCGGCTGC	CAGCTGCTGC
1021	TACTGA			

156/156

//

## INTERNATIONAL SEARCH REPORT

International application No.  
PCT/US93/10507

**A. CLASSIFICATION OF SUBJECT MATTER**

IPC(5) : G06F 15/42

US CL : 364/413.01

According to International Patent Classification (IPC) or to both national classification and IPC

**B. FIELDS SEARCHED**

Minimum documentation searched (classification system followed by classification symbols)

U.S. : 364/413.01; 435/6; 536/23.1

Documentation searched other than minimum documentation to the extent that such documents are included in the fields searched

Electronic data base consulted during the international search (name of data base and, where practicable, search terms used)  
MEDLINE, CA

**C. DOCUMENTS CONSIDERED TO BE RELEVANT**

Category*	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.
P,X	GENETIC ENGINEERING NEWS, Vol. 13, No. 18, issued 15 October 1993, Potera, "Hitachi Chemical Offers Probe Design Software and Service," pages 1, 22, see entire document.	1-102
Y	NUCLEIC ACIDS RESEARCH, Vol. 18, No. 7, issued 11 April 1990, Lowe et al., "A Computer Program for Selection of Oligonucleotide Primers for Polymerase Chain Reactions," pages 1757-1761, see entire document.	1-102
Y	METHODS IN ENZYMOLOGY, Vol. 183, issued 1990, Landau et al., "Fast Alignment of DNA and Protein Sequences," pages 487-502, see entire document.	1-102

☐ Further documents are listed in the continuation of Box C.

☐ See patent family annex.

* Special categories of cited documents:	* T	later document published after the international filing date or priority date and not in conflict with the application but cited to understand the principle or theory underlying the invention
* A* document defining the general state of the art which is not considered to be part of particular relevance	* X*	document of particular relevance; the claimed invention cannot be considered novel or cannot be considered to involve an inventive step when the document is taken alone
* E* earlier documents published on or after the international filing date	* Y*	document of particular relevance; the claimed invention cannot be considered to involve an inventive step when the document is combined with one or more other such documents, such combination being obvious to a person skilled in the art
* L* document which may throw doubts on priority claim(s) or which is cited to establish the publication date of another citation or other special reason (as specified)	* &*	document member of the same patent family
* O* document referring to an oral disclosure, use, exhibition or other means		
* P* document published prior to the international filing date but later than the priority date claimed		

Date of the actual completion of the international search

21 December 1993

Date of mailing of the international search report

03 FEB 1994

Name and mailing address of the ISA/US  
Commissioner of Patents and Trademarks  
Box PCT  
Washington, D.C. 20231

Authorized officer

SCOTT HOUTTEMAN

Facsimile No. NOT APPLICABLE

Telephone No. (703) 308-0196

Form PCT/ISA/210 (second sheet)(July 1992)\*

CGK00009339